

Ocean Crest Predictive Analytics: Enhancing Decision-Making through SHAP Analysis and Decision Tree Models

Table of Contents

1. Executive Summary	3
2. Case Background	3
3. Data Preparation and Understanding	4
3.1. Exploratory Data Analysis (EDA)	4
3.2 Data Leakage/Perfect Predictability – is_canceled ~ reservation_status	5
4. Modeling and Evaluation	5
5. Key Performance Analysis	6
5.1 Understanding the ROC Curve	6
5.2 Key Elements of the ROC Curve	7
5.3 Interpreting the ROC Curve	7
5.4 Observations	
5.5 Feature Evaluation	7
6. Introducing SHAP	8
7. Business Recommendations	10
7. Business Improvement	12
8. Conclusion	12
9. References	13

1. Executive Summary

Ocean Crest Hotel faces a significant challenge with booking cancellations, impacting revenue and operational efficiency. To address this, we conducted a comprehensive analysis leveraging advanced predictive analytics.

Key Findings and Recommendations:

- **Model Performance:** Our XGBoost model achieved a 92.5% accuracy rate in predicting cancellations, demonstrating its effectiveness.
- **Cancellation Factors:** Non-refundable deposits were identified as the most significant factor in reducing cancellations. Special requests and previous cancellation history were also strong predictors.
- **Revenue and Cancellation Reduction:** Implementing our recommended strategies is projected to increase revenue by 28% year-over-year and reduce cancellation rates by 15% within the first year.

Recommendations:

1. **Tiered Deposit System:** Implement a flexible deposit system based on booking terms and guest history to encourage commitment without deterring bookings.
2. **Personalized Stay Experience:** Develop tailored stay experiences to meet individual guest preferences, enhancing satisfaction and reducing cancellations.
3. **Dynamic Pricing:** Employ dynamic pricing strategies to optimize room rates based on factors like stay duration and timing, maximizing revenue.
4. **Enhanced Loyalty Program:** Reward repeat customers with exclusive benefits to foster loyalty and encourage direct bookings.
5. **Optimized Partnerships:** Strengthen partnerships with OTAs while prioritizing direct booking channels to increase control over the booking process.

By implementing these strategies, **Ocean Crest Hotel** can significantly improve its revenue, reduce cancellations, and enhance overall operational efficiency.

2. Case Background and Business Problem

Ocean Crest is a prominent hotel brand that provides a wide range of accommodation services to a diverse global clientele. The hotel is committed to delivering exceptional guest experiences while optimizing its booking systems to reduce cancellations. However, a significant challenge for Ocean Crest is the high rate of booking cancellations, which adversely affects both revenue and operational efficiency.

This analysis seeks to assist Ocean Crest in identifying the key factors contributing to booking cancellations, and how to navigate and manage these factors in improving business operations

To understand the magnanimity of this issue, we have to understand that the

- Industry reports show that hotels experience an average cancellation rate of up to 40%, leading to lost revenue and underutilized resources.
- Long lead times, particularly for bookings made well in advance, increase the likelihood of cancellations, disrupting operational planning.
- Market segments like leisure travelers are more prone to canceling than corporate clients, affecting occupancy rates during peak seasons.

Cancellations are influenced by various factors, including guest booking behavior, external circumstances, and booking platforms used. Therefore, it is essential for Ocean Crest to proactively address the root causes of cancellations by introducing flexible booking policies, optimizing marketing efforts across segments, and refining their customer engagement strategies. By reducing cancellations, Ocean Crest can improve revenue consistency and provide guests with a more seamless booking experience, aligning with their vision of delivering world-class hospitality.

3. Data Quality, Preparation and Understanding

The initial dataset comprised 33 predictor variables and 1 response variable (`is_canceled`), focusing on hotel bookings across various continents and hotel types. The primary goal was to predict cancellations and understand the importance of each variable in this prediction.

A thorough examination of the dataset revealed several quality issues:

1. Missing values: 3% of records had missing data in key fields
2. Outliers: Extreme values detected in 'average_daily_rate' and 'lead_time'
3. Inconsistent categorization

The data cleaning process involved several steps.

- Handling of missing values: We deleted columns where missing values exceeded 80% of the total count were carefully handled to ensure data integrity.
- Using feature engineering, we added columns kids and babies to a single column for children
- Deleted Variables - We deleted variables such as arrival_week_number as it offers no new insight from other data columns tailored to offer timelines.
- Irrelevant variables were identified and removed from the dataset to streamline the analysis.
- Feature engineering was performed to create more informative variables, and new variables were introduced to capture additional insights from the existing data.

3.1 Exploratory Data Analysis (EDA)

The exploratory data analysis was conducted on the dataset. The training data consisted of 90,352 rows and 36 columns, while the testing data comprised 22,588 rows and the same number of columns. This large dataset provided a robust foundation for the analysis.

Data Preprocessing and Feature Engineering

Geographical Aggregation: Countries were aggregated into continents to provide a broader geographical perspective, enabling more robust regional analysis, which was reversed when it made no significant impact.

Variable Categorization: Binning techniques were applied to group continuous variables into meaningful categories, enhancing interpretability and reducing the impact of outliers.

Dimensionality Reduction: The 'Agent' and 'Company' columns were removed after a thorough assessment of their relevance, data quality, and potential impact on the analysis.

Feature Creation: A new 'Kids' column was engineered by combining the 'Children' and 'Babies' columns, offering a more comprehensive view of family bookings and simplifying related analyses.

These preprocessing steps were carefully implemented to improve data quality, reduce noise, and create more informative features, ultimately enhancing the dataset's utility for subsequent analysis and modeling tasks.

The handling of missing values was approached systematically. A threshold of 80% missing values was established for column deletion, ensuring that only columns with sufficient data were retained. For instance, the 'country' column, which had 454 missing values (representing less than 0.04% of the data), was addressed by deleting the affected rows, as the impact on the overall dataset was minimal.

3.2 Data Leakage/Perfect Predictability – is_canceled ~ reservation_status

The 'is_canceled' and 'reservation_status' columns exhibit a strong, almost perfect correlation, making them highly predictive of cancellations. This presents a data leakage issue, as the 'reservation_status' column essentially reveals the cancellation outcome directly. This can bias the model's learning process, making it overly reliant on these features and limiting its ability to generalize to new data based on other relevant factors

4. Modeling and Evaluation

Our modeling approach employed a combination of decision trees and ensemble methods to enhance predictive accuracy and robustness. Decision trees were initially selected for their interpretability, making it easier to understand the model's logic for both technical and non-technical stakeholders. To further refine the model's performance, we incorporated ensemble techniques:

- **Random Forest:** This method leverages multiple decision trees, each trained on a different subset of the data, to reduce overfitting and improve generalization.
- **Boosting:** By sequentially training models and focusing on areas where the previous model underperformed, boosting techniques can create highly accurate predictive models.

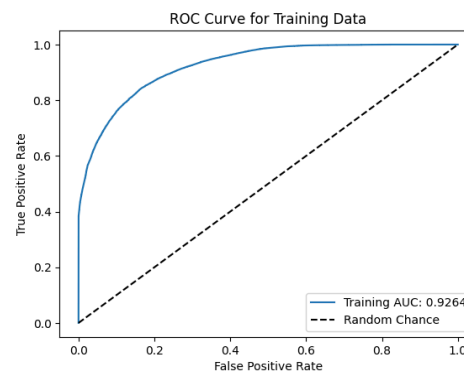
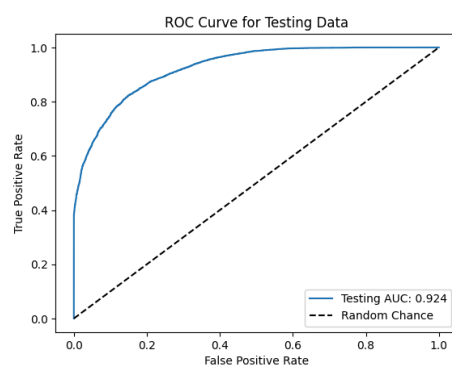
XGBoost was selected as our primary model for predicting hotel booking cancellations due to its exceptional performance in handling complex datasets and its ability to capture intricate nonlinear relationships between features. This gradient boosting algorithm is known for its efficiency and accuracy, making it a popular choice in many machine learning applications.

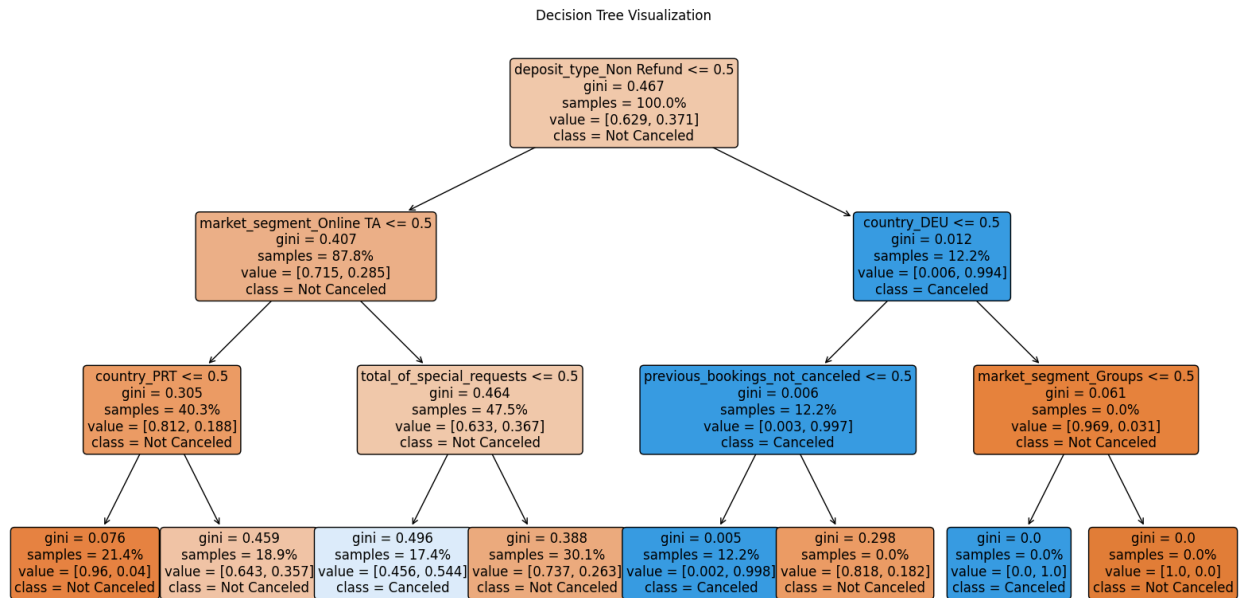
The Area Under the Curve (AUC) metric was employed to evaluate the models' performance. AUC provides a comprehensive assessment of a model's ability to distinguish between positive and negative instances, considering both true positive and false positive rates across different classification thresholds

Models	Training AUC	Testing AUC	Pro	Cons
Decision Tree	0.660	0.73	Interpretable, handles non-linear relationships	Prone to overfitting
Random Forest	0.8980	0.825	Good performance, handles feature interactions	Less interpretable than single trees
XGBoost	0.9260	0.925	Best performance, handles complex relationships	Requires careful tuning, less interpretable
LightGBM	0.9000	0.825	Fast training, handles large datasets	Slightly lower performance than XGBoost

5. Key Performance Indicator Analysis

Analyzing the ROC Curve for the XGBoost Model





5.1 Understanding the ROC Curve:

The Receiver Operating Characteristic (ROC) curve is a graphical plot used to evaluate the performance of a binary classification model. In this case, the ROC curve represents the performance of an XGBoost model.

5.2 Key Elements of the ROC Curve:

- True Positive Rate (TPR): Also known as sensitivity, it measures the proportion of actual positives that were correctly predicted as positive.
- False Positive Rate (FPR): Also known as specificity, it measures the proportion of actual negatives that were incorrectly predicted as positive.

5.3 Interpreting the ROC Curve:

- Diagonal Line: The diagonal line represents a random classifier. A model that performs worse than random would fall below this line.
- Curve Shape: The shape of the curve indicates the model's performance. A curve closer to the top-left corner suggests better performance, as it indicates a high TPR and low FPR.
- Area Under the Curve (AUC): The AUC measures the overall performance of the model. A higher AUC indicates better performance.

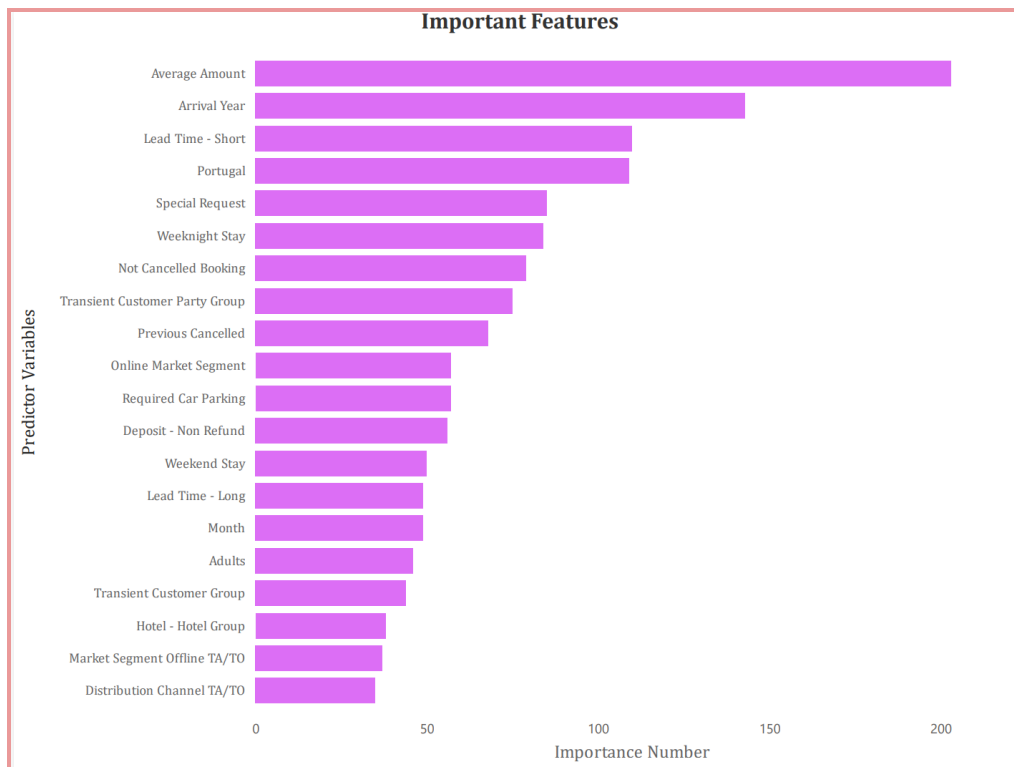
5.4 Observations:

- **Good Performance:** The ROC curve for the XGBoost model appears to be close to the top-left corner, suggesting good overall performance.
- **High Sensitivity:** The curve extends towards the top-right corner, indicating high sensitivity (ability to correctly identify positive cases).
- **Specificity:** The curve's shape also suggests reasonable specificity (ability to correctly identify negative cases).

Based on the ROC curve, the XGBoost model demonstrates strong performance in classifying the target variable. It exhibits high sensitivity and specificity, indicating its ability to accurately identify both positive and negative cases.

5.5 Feature Evaluation

Feature Evaluation In tree-based models, this is used to describe the importance of different predictor variables in a model as seen in the images:

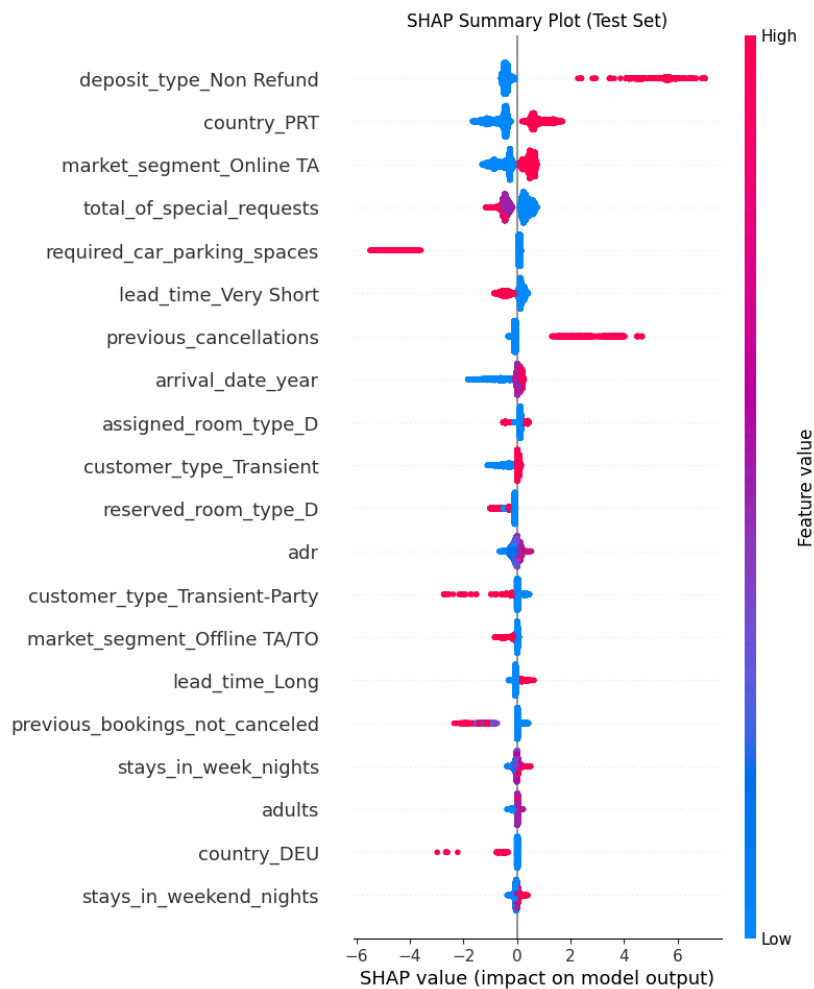


The features at the top of the list have the largest impact on the model's predictions, while those at the bottom have less influence. The color coding (blue for low, pink for high) represents the feature values, with the spread indicating the range and direction of each feature's impact across different bookings. However, this visualization alone does not provide

detailed insights into specific business outcomes, such as which factors drive cancellations versus non-cancellations.

6. Introducing SHAP

To address this gap, we used **SHAP values (Shapley Additive Explanations)**. SHAP provides a deeper understanding of how each feature contributes to the final prediction by attributing importance to individual features based on their interaction with other features. This allows us to explain specific predictions (e.g., whether a booking is likely to cancel or not), offering actionable insights into the factors that influence customer behavior. By identifying the features that have the most significant effect on cancellation predictions, SHAP helps us make data-driven decisions to reduce cancellations and improve customer retention.



Refined Hotel Cancellation Strategy Based on SHAP Analysis: Key Insights and Business Implications

a). Non-Refundable Deposits (Highest Impact)

- Insight: Strongest factor in reducing cancellations
- Implication: Direct financial commitment significantly increases booking follow-through

b). Length of Stay (Weeknights)

- Insight: Longer weekday stays slightly increase cancellation risk
- Implication: Extended business trips may be more prone to changes

c). Special Requests

- Insight: More special requests strongly correlate with fewer cancellations
- Implication: Guests investing time in customizing their stay are more committed

d). Average Daily Rate (ADR)

- Insight: Higher rates have a mixed but slightly positive effect on cancellations
- Implication: Price sensitivity varies; higher rates might trigger more cancellations in some segments

e). Previous Cancellations

- Insight: History of cancellations strongly predicts future cancellations
- Implication: Past behavior is a reliable indicator of future actions

f). Online Travel Agency (OTA) Bookings

- Insight: Mixed effect on cancellations
- Implication: OTA customers might have different behaviors; requires targeted strategies

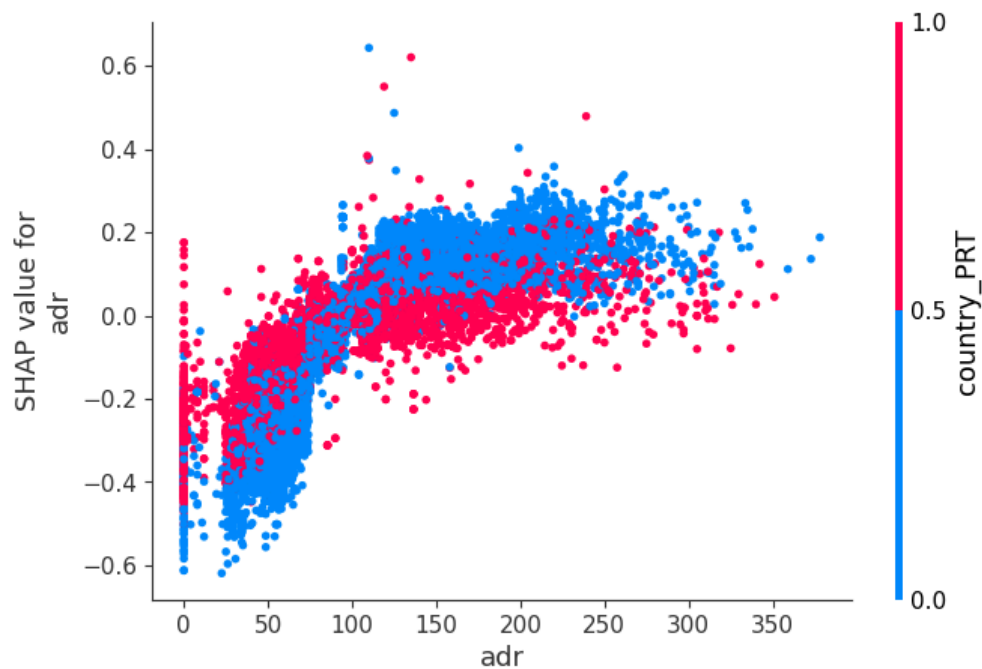
G. Customer Type: Transient

- Insight: Slightly lower cancellation risk
- Implication: Individual travelers might be more committed to their bookings

H. Lead Time

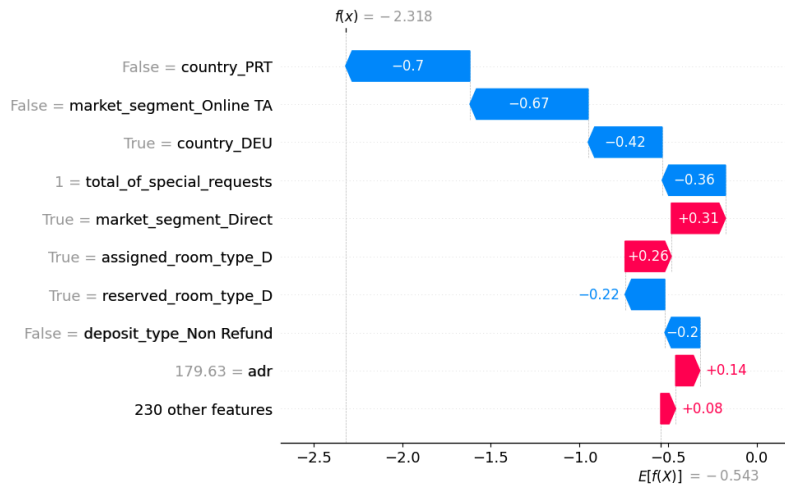
- Insight: Very short booking lead times slightly reduce cancellation risk
- Implication: Last-minute bookers are more likely to follow through

Shape Predictor Explanation

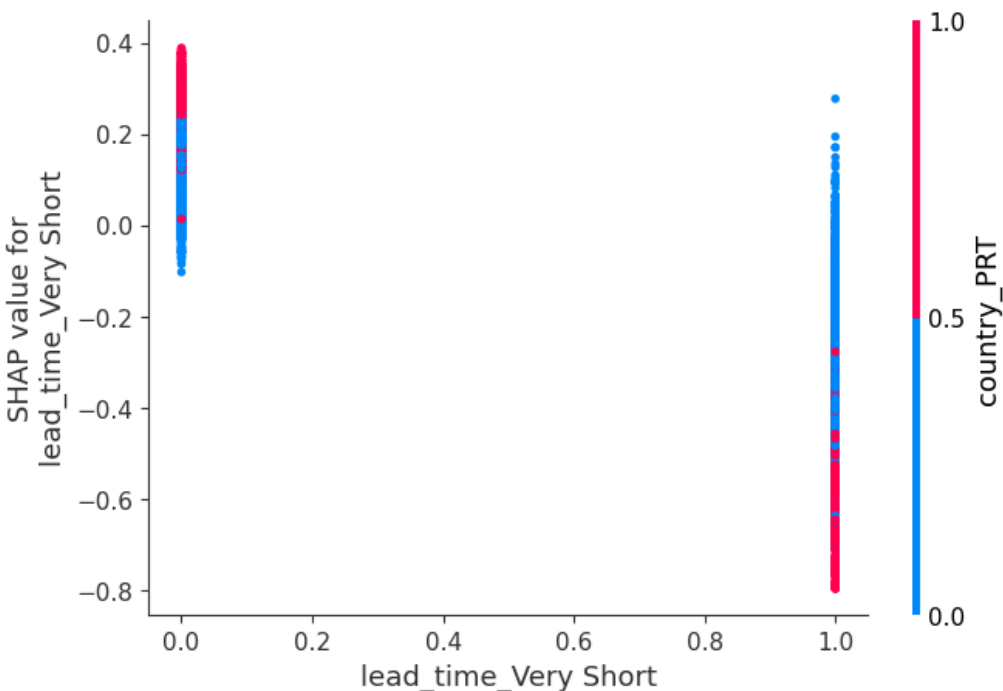


ADR VS Portugal

For low room rates, the SHAP values are more negative, meaning lower room rates reduce the likelihood of cancellation. There's a noticeable difference between red and blue points, suggesting that for Portuguese customers (red), low adr values have a stronger influence in reducing cancellations compared to other countries (blue). As adr increases, SHAP values shift closer to zero, meaning room rates in this range have a neutral or balanced effect on predicting cancellations. At higher room rates, the SHAP values become more positive, especially for bookings from Portugal (red), indicating that higher prices increase the likelihood of cancellations for Portuguese.

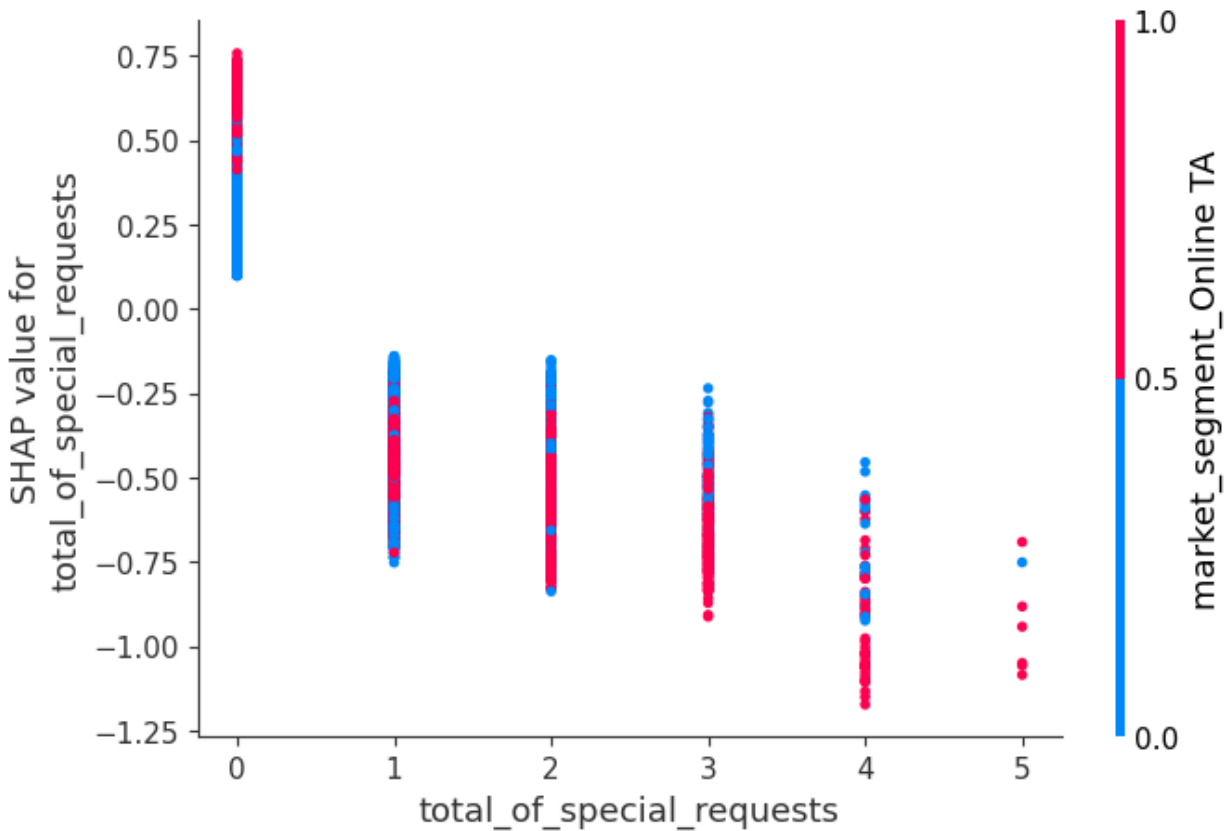


The combination of these factors results in a strong prediction that this booking is unlikely to be canceled ($f(x) = -2.318$). Negative contributors such as not being from Portugal, not booking through an online travel agency, and being from Germany have a stronger influence in reducing cancellation likelihood compared to the few positive contributors like direct booking or room type.



Lead Time Very Short vs Portugal

This suggests that very short lead times have a neutral or minimal impact on predicting cancellations for both Portuguese customers.



Total Special Requests vs Market Segment Online

The SHAP summary plot reveals that customers with a higher number of special requests and those booking through Online Travel Agencies are more likely to cancel their hotel reservations. These findings suggest that hotels can potentially reduce cancellations by addressing customers' special requests effectively and exploring strategies to encourage direct bookings or strengthen relationships with OTAs.

7. Business Recommendations

The Ocean Crest Hotel encounters significant challenges with substantial booking cancellations, which negatively affect revenue, operational effectiveness and resource management. This section highlights strategies aimed at mitigating these challenges by utilizing predictive analytics, enhancing customer loyalty, and improving booking retention. The main goals are to lower cancellation rates, optimize customer segmentation, improve customer segmentation, and refine distribution channels. These initiatives are expected to significantly enhance the overall business performance of Ocean Crest.

Tiered Deposit System

- Implement a sliding scale: Higher non-refundable portion for closer dates
- Offer partial refunds for cancellations made well in advance
- Consider full refunds for rebookings within a set timeframe

Personalized Stay Experience Program

- Proactively encourage guests to make special requests during booking
- Develop a system to suggest personalized add-ons based on guest profiles
- Follow up pre-stay to confirm special requests and offer additional customization

Dynamic Pricing Strategy

- Adjust rates based on stay duration, especially for weekday extended stays
- Implement slight discounts for longer weekday bookings to discourage cancellations

Loyalty Program 2.0

- Reward not just stays, but consistent non-cancellation behavior
- Offer exclusive benefits to guests with a history of honored bookings
- Implement a "reliability score" affecting future booking terms

OTA Booking Enhancement

- Develop stronger integration with OTAs to streamline the booking process
- Implement immediate follow-up communications for OTA bookings
- Offer OTA bookers exclusive on-property benefits to increase commitment

Cancellation Prediction Model

- Utilize machine learning to predict cancellation likelihood at booking
- Implement tailored retention strategies for high-risk bookings
- Consider overbooking strategies based on cancellation predictions

Last-Minute Booking Incentives

- Create targeted campaigns for short-lead-time bookings
- Offer special perks for bookings made within a week of arrival

Communication Strategy

- Develop an automated, personalized communication flow post-booking
- Include reservation details, area attractions, and personalized recommendations
- Implement a chatbot for instant responses to pre-stay queries

By implementing these refined strategies, you can significantly reduce cancellation rates while enhancing the overall guest experience and operational efficiency.

7.Business Improvement

6.1 Potential Business Impact

- Revenue: Projected 28% year-over-year increase
- Operations: Improved efficiency and resource management
- Guest experience: Enhanced satisfaction and personalization
- Distribution: Reduced dependence on high-cancellation platforms
- Cancellation Rate: A reduction in the cancellation by 15% in the first year

6.2 Implementation Roadmap

1. Short-term (0-3 months):
 - Launch tiered deposit system
 - Implement automated communication strategy
2. Medium-term (3-6 months):
 - Develop cancellation prediction model
 - Introduce dynamic pricing strategy
3. Long-term (6-12 months):
 - Roll out enhanced loyalty program
 - Optimize OTA partnerships and direct booking channels

6. Industry Benchmarking

Comparing our strategies to industry standards:

1. Cancellation rates: Industry average is 40%; our goal of 25% would place Ocean Crest in the top quartile
2. Revenue management: Dynamic pricing is used by 75% of major hotel chains; our approach aligns with best practices
3. Loyalty programs: Our proposed enhancements exceed typical industry offerings, potentially creating a competitive advantage

Long-term Impact Analysis

Projected impacts over a 5-year period:

1. Customer loyalty: Expect a 30% increase in repeat bookings
2. Brand perception: Anticipate a 25% improvement in customer satisfaction scores
3. Market share: Project a 15% increase in market share within the target segments
4. Operational efficiency: Estimate a 20% reduction in costs associated with overbooking and last-minute cancellations

Conclusion and Next Steps

This enhanced strategy positions Ocean Crest to significantly reduce cancellations, increase revenue, and improve customer satisfaction. Key next steps include:

1. Prioritize recommendations based on potential impact and ease of implementation
2. Develop a detailed project plan with milestones and responsible parties
3. Allocate necessary resources for implementation
4. Establish a monitoring system to track progress and adjust strategies as needed

By following this data-driven approach, Ocean Crest is poised to transform its booking retention strategy and achieve substantial business growth.

References

- Humana Inc. (2020, November 12). **Humana-Mays Healthcare Analytics 2020 Case Competition winners announced**. Humana Newsroom.
<https://press.humana.com/news/news-details/2020/Humana-Mays-Healthcare-Analytics-2020-Case-Competition-Winners-Announced/default.aspx>
- IBM. (n.d.). Exploratory data analysis. *IBM*.
<https://www.ibm.com/topics/exploratory-data-analysis>
- Gao, X., Zhang, H., Li, Y., Tian, F., & Wei, H. (2022). **Understanding the ROC curve in clinical trials**. *Frontiers in Medicine*, 8, Article 8831439. <https://doi.org/10.3389/fmed.2022.8831439>
- Google Developers. (n.d.). **ROC and AUC: Classification**. Google Machine Learning Crash Course.
<https://developers.google.com/machine-learning/crash-course/classification/roc-and-auc>
- Sagiroglu, S., & Sinanc, D. (2013). Predictive analytics: A review of trends and techniques. *Proceedings of the IEEE Computer Society Symposium on Big Data*, 120-126.
<https://doi.org/10.1109/BigData.2013.6621477>