# UNIVERSITY OF TEXAS AT ARLINGTON
# INSY-5339-003-PRIN OF BUSINESS DATA MINING
# SPRING 2022

# FINAL PROJECT REPORT

**Due Date: 5/8/2022**

**Submitted by: Team 9**

**Team Members**
Ruchi Shukla: 1001977095
Pranati Chauhan: 1002032137
Jessica Mendem: 1002036921

**INDEX**

# Executive Summary

**Company Goal:**

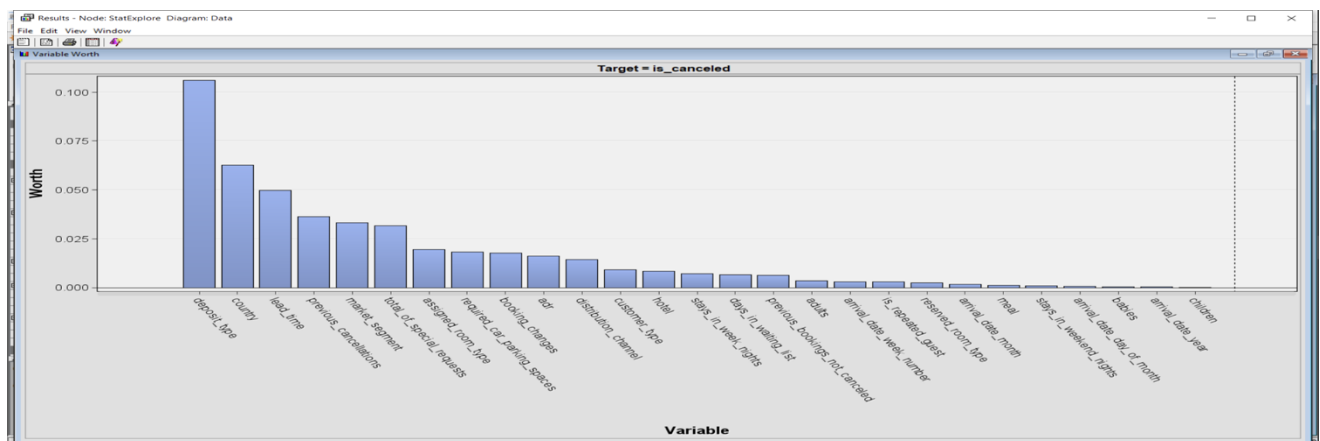Boost business for hotels by predicting cancellations.

**Problem Statement:**

In the hospitality industry, there is a vast amount of data which is untouched that can revolutionize revenue management. In the industry, optimizing the operations of the hotel industry to ensure the right room to the right customer at the right price is what would hotels stay ahead of the competition and help in efficient use of resources.

Last-minute cancellations and no-show lead to heavy loss of resources. What if cancellations could be predicted way before they are made?

**Proposed Solution:**

If potential cancellers are identified, they may be persuaded to check into a hotel by attractive, customized incentives. Analyzing the behavior of previous hotel guests might help management provide targeted incentives, allowing them to keep as many reservations as feasible.



As seen in the bar graph, it can be concluded that deposit, country, and lead time are the top 3 factors which lead to cancellations.

A proper solution investigating these factors can help to minimize cancellations.

**Approach:**

To implement the concept, data was analyzed from three years (2015-2017) of customer records at city and resort hotels. To construct the best accurate prediction model for identifying probable cancellers and studying prior customer behavior, a variety of data mining algorithms were applied.

**Results:**

To detect prospective cancellers, a predictive model with a high degree of accuracy was successfully created. Various parameters were used to assess consumer behavior and create client categories. Incentives were given to customers who were likely to cancel reservations.

**Managerial Insights:**

Managers might utilize the prediction model to target possible cancellers, and group of customers analysis could be used to give individualized rewards to the identified potential cancellers. The benefits must be based on the resources and capital available.

# Hotel Booking Cancellations

**Background:**

We are a data analytics team for a hotel chain, and our job is to provide management with insights that will help them enhance business at their city and resort properties.

• Numerous individuals who reserve hotel reservations end up checking in, but there are also many people who cancel their reservations.

• It is wonderful for business if people show up; but it is a loss for the franchise if they cancel their reservations and do not show up.

**Solution:**

Management has advised us that if we can identify possible cancellers and provide them with individualized incentives, they may be persuaded to visit a hotel, minimizing the chance of a reservation cancellation.

We'll discover and investigate folks who have a high risk of canceling in order to provide them with targeted incentives.

**What Should You Do**

- If we can identify potential cancellers, we might be able to convince them and provide specific incentives, minimizing the risk of a reservation cancellation.

- We'll discover and investigate folks who have a high risk of canceling to provide them with targeted incentives.

# Description of Data

**Resources on hand**

Data was collected from Kaggle.

•This dataset contains client records of 2 hotels- city hotel and a resort hotel (2015-2017).

•This dataset comprises booking information including dates of booking, duration of stay, number of adults, children, and/or newborns, meal preferences, room type, etc.

**List of Variables**

| Variable | Description | Type |
|---|---|---|
| Is_canceled | If booking was cancelled | Categorical |
| lead_time | difference between booking date & arrival date | Integer |
| arrival_date_year | arrival date of the year | Integer |
| arrival_date_month | arrival month - Jan to Dec | Categorical |
| arrival_date_week_number | week of arrival | Integer |
| arrival_date_day_of_month | date and day of arrival | Integer |
| stays_in_weekend_nights | number of nights stayed on weekends | Integer |
| stays_in_week_nights | number of nights stayed in a week | Integer |
| Adults | Number of adults | Integer |
| Children | Number of children | Integer |
| Babies | Number of babies | Integer |

| Meal | type of meal booked | Categorical |
|------|---------------------|-------------|
| Country | country of origin of the guest | Categorical |
| market_segment | source of booking - self or through travel agent or other sources | Categorical |
| distribution_channel | source of booking - self or through travel agent or other sources | Categorical |
| is_repeated_guest | If the guest | Categorical |
| previous_cancellations | bookings cancelled by the customer at the same hotel before the current booking | Integer |
| previous_bookings_not_canceled | number of bookings not cancelled by guest, prior to current booking | Integer |
| reserved_room_type | which type of room was booked by the customer | Categorical |
| assigned_room_type | which room was assigned | Categorical |
| booking_changes | count of changes made to the reservation ahead of check-in | Integer |
| deposit_type | Category addressing whether a deposit was made during reservation or not | Categorical |
| Agent | travel agency identifier/code | Categorical |
| Company | company identifier/code | Categorical |
| days_in_waiting_list | measurement of time when booking was on waitlist | Integer |
| customer_type | type of customer | Categorical |

| | | |
|---|---|---|
| Adr | Average daily rates at the time of booking | Numeric |
| required_car_parking_spaces | number of parking spaces requested in the booking | Integer |
| total_of_special_requests | count of special requests associated with the reservation | Integer |
| reservation_status | status of reservation at the end of stay | Categorical |
| reservation_status_date | date of reservation_status | Date |

**Summary of Variables**

| | |
|---|---|
| Total number of records | 118,221 |
| Binary variable | 2 |
| Ordinal variable | 4 |
| Nominal | 9 |
| Interval | 13 |
| Outcome/Target variable | Is_canceled |
| Percentage of the records that belong to each class. | For is_canceled = 0(Not Canceled):74283 Observations = 63% For is_canceled = 1(Canceled): 43938 Observations = 37% |

# Data Visualization

For data exploration and finding trends and patterns, data visualization is a necessary tool. It helps us understand relationships between variables.

Through these visualizations, univariate analysis and various patterns and trends helped in understanding and brainstorming why cancellations happen and how changes can be made by analyzing the trends.

## Visualization Methods:

- Bar Chart
- Pie Chart
- Heat Map

## Univariate Analysis



The univariate analysis here shows the following:

- City Hotel had more reservations and cancellations than Resort Hotels.
- Around 49% of the bookings were made by Transient customer type and had a 69% rate of cancellations in city hotels.
- Resort Hotels had most bookings and cancellations by Transient Customers as well.

## Maps

Sheet 2

Top 10 Countries with Bookings

- Portugal has the most bookings and cancellations.
- Among 47.887% of bookings, 56.9% were cancelled, which is an alarming percentage.

## Bar Graph



Booking % Per Year

- **Booking percentage per year:** In 2015 bookings are 62.8% and Cancellation 37.2%. The highest bookings were in 2016 bookings are 64.1% and cancellations are 35.9%. In 2017 bookings were up to 61.1% and 38.9% Cancellations.

Busiest Month

- **Busiest Month:** In overall months we have busiest month is in **August** and it also has the highest cancellation rate.



Reserved Room Type

- Customers who booked **Room A** cancelled more than any other room.

## Heat Map



- It was found that the variables were weakly correlated with each other.

# Data Mining Models

The data was utilized to create a model that predicted how likely a potential customer would cancel their reservation.

For developing the prediction model:

- Naïve Bayes Classification
- K- Nearest Neighbor Algorithm
- Classification Tree
- Logistic Regression

- After partitioning the dataset into a 60:20:20 ratio (Training: Validation: Test), all the models were run.
- **Naive Bayes and K-Nearest Neighbor** algorithms were used, because our dataset was vast, and these algorithms perform better on large datasets.
- In case the above-mentioned models' performance was unsatisfactory, **classification tree and logistic regression** were used as backups.
- We also performed different regression models- stepwise, forward, and backward, and the misclassification rate for forward regression was minimum.

**Process Flow Diagram**

# Model 1: Naïve Bayes Classification:

## Results



## Statistics and Bayesian Network

## Observations

- Average Squared Error for validation dataset **= 17%**
- Misclassification Rate for validation dataset **= 24%**

## Important Variables in Bayesian Network:

- lead_time
- market_segment
- assigned_room_type
- customer_type
- required_car_parking_spaces
- country
- hotel
- is_repeated_guest
- meal
- reserved_room_type
- deposit_type
- distribution_channel
- total_of_special_requests

**Model 2: KNN Algorithm -** Default settings were used.

**Results**

## Confusion Matrix



## Observations

- Average Squared Error for validation dataset **= 18%**
- Misclassification Rate for validation dataset **= 26%**

## Model 3: **Classification Tree-** Default settings were used.

## Results



## Observations:

- Average Squared Error for validation dataset **= 13%**
- Misclassification Rate for validation dataset **= 20%**

# List of important variables

Variable Importance

| Variable Name | Label | Number of Splitting Rules | Importance | Validation Importance | Ratio of Validation to Training Importance |
|---|---|---|---|---|---|
| deposit_type | | 1 | 1.0000 | 1.0000 | 1.0000 |
| country | | 2 | 0.4782 | 0.4715 | 0.9861 |
| market_segment | | 1 | 0.4289 | 0.4214 | 0.9824 |
| lead_time | | 3 | 0.3975 | 0.4041 | 1.0166 |
| total_of_special_requests | | 1 | 0.3444 | 0.3642 | 1.0577 |
| previous_cancellations | | 3 | 0.2811 | 0.3006 | 1.0695 |
| required_car_parking_spaces | | 1 | 0.2574 | 0.2625 | 1.0200 |
| booking_changes | | 1 | 0.1868 | 0.2015 | 1.0787 |
| arrival_date_year | | 1 | 0.1484 | 0.1412 | 0.9516 |
| previous_bookings_not_canceled | | 2 | 0.0786 | 0.0858 | 1.0925 |
| distribution_channel | | 1 | 0.0710 | 0.0608 | 0.8556 |
| customer_type | | 1 | 0.0513 | 0.0537 | 1.0461 |
| hotel | | 1 | 0.0337 | 0.0292 | 0.8658 |

## Model 4: Logistic Regression – Default settings were used.

## Results

Fit Statistics

Target=is_canceled Target Label=' '

| Fit Statistics | Statistics Label | Train | Validation | Test |
|---|---|---|---|---|
| _AIC_ | Akaike's Information Criterion | 53505.06 | . | . |
| _ASE_ | Average Squared Error | 0.12 | 0.12 | 0.12 |
| _AVERR_ | Average Error Function | 0.37 | 0.39 | 0.38 |
| _DFE_ | Degrees of Freedom for Error | 70696.00 | . | . |
| _DFM_ | Model Degrees of Freedom | 235.00 | . | . |
| _DFT_ | Total Degrees of Freedom | 70931.00 | . | . |
| _DIV_ | Divisor for ASE | 141862.00 | 47290.00 | 47290.00 |
| _ERR_ | Error Function | 53035.06 | 18289.20 | 18145.51 |
| _FPE_ | Final Prediction Error | 0.12 | . | . |
| _MAX_ | Maximum Absolute Error | 1.00 | 1.00 | 1.00 |
| _MSE_ | Mean Square Error | 0.12 | 0.12 | 0.12 |
| _NOBS_ | Sum of Frequencies | 70931.00 | 23645.00 | 23645.00 |
| _NW_ | Number of Estimate Weights | 235.00 | . | . |
| _RASE_ | Root Average Sum of Squares | 0.35 | 0.35 | 0.35 |
| _RFPE_ | Root Final Prediction Error | 0.35 | . | . |
| _RMSE_ | Root Mean Squared Error | 0.35 | 0.35 | 0.35 |
| _SBC_ | Schwarz's Bayesian Criterion | 55659.89 | . | . |
| _SSE_ | Sum of Squared Errors | 17216.67 | 5819.42 | 5870.64 |
| _SUMW_ | Sum of Case Weights Times Freq | 141862.00 | 47290.00 | 47290.00 |
| _MISC_ | Misclassification Rate | 0.18 | 0.18 | 0.18 |

Classification Table

Data Role=TRAIN Target Variable=is_canceled Target Label=' '

| Target | Outcome | Target Percentage | Outcome Percentage | Frequency Count | Total Percentage |
|---|---|---|---|---|---|
| 0 | 0 | 82.1738 | 90.7694 | 40455 | 57.0343 |

Type 3 Analysis of Effects

| Effect | DF | Wald Chi-Square | Pr > ChiSq |
|---|---|---|---|
| IMP_children | 1 | 35.4431 | <.0001 |
| adr | 1 | 346.1667 | <.0001 |
| adults | 1 | 49.1659 | <.0001 |
| arrival_date_day_of_month | 1 | 57.7395 | <.0001 |
| arrival_date_month | 11 | 293.6240 | <.0001 |
| arrival_date_week_number | 1 | 60.3100 | <.0001 |
| arrival_date_year | 0 | 0.0000 | . |
| assigned_room_type | 11 | 922.4156 | <.0001 |
| babies | 1 | 0.5301 | 0.4666 |
| booking_changes | 1 | 262.8020 | <.0001 |
| country | 168 | 5314.6726 | <.0001 |
| customer_type | 3 | 388.2299 | <.0001 |
| days_in_waiting_list | 1 | 13.1290 | 0.0003 |
| deposit_type | 2 | 989.6810 | <.0001 |
| distribution_channel | 3 | 717.1956 | <.0001 |
| hotel | 1 | 5.2996 | 0.0213 |
| is_repeated_guest | 1 | 104.8284 | <.0001 |
| lead_time | 1 | 1730.1400 | <.0001 |
| market_segment | 7 | 1665.8151 | <.0001 |
| meal | 3 | 92.7265 | <.0001 |
| previous_bookings_not_canceled | 1 | 164.7710 | <.0001 |
| previous_cancellations | 1 | 706.1083 | <.0001 |
| required_car_parking_spaces | 1 | 17.7431 | <.0001 |
| reserved_room_type | 8 | 631.7160 | <.0001 |
| stays_in_week_nights | 1 | 90.9505 | <.0001 |
| stays_in_weekend_nights | 1 | 64.6551 | <.0001 |
| total_of_special_requests | 1 | 2205.5688 | <.0001 |

# Comparison between different Regression



# Confusion Matrix



# Observations

- Average Squared Error for validation dataset **= 12%**
- Misclassification Rate for validation dataset **= 18%**

## Model Comparison

- K-NN model was predicting incorrectly.
- Logistic regression and classification tree were chosen since they performed much better.
- Misclassification rate for the logistic regression model was **18%.**



**Clustering -** 3 significant clusters were formed: Cluster 1, 2, 3

Cluster 2: Customers likely to cancel (1)

Cluster 3: Customers not likely to cancel (0)

| is_canceled=0 | is_canceled=1 |
|---|---|
| 0.540964 | 0.459036 |
| 0.351651 | 0.648349 |
| 0.725227 | 0.274773 |



**Segment 2: 65%** consists of people who canceled reservations.

**Statistics for variables**:

- Lead time – 329
- Market segment – Groups = 48%, Offline TA/TO = 31%
- Deposit Type – No Deposit = 59%
- Average Daily Rate = $84 .00
- Late check ins
- Highest number of previous cancellations
- Travelers were mostly from **Portugal** (59%) and **Germany** (10%)
- Most arrivals in **August** and **September**
- Waiting list was higher
- Reserved Room Type A – 88%
- People who booked Bed & Breakfast – 77%
- Booked City Hotel – 74%
- Transient customers – 54%
- Distribution Channel – TA/TO = 93%

# Interpretation and Results

- To identify potential cancellations, logistic regression or classification tree can be used.
- Customers can be segmented to offer specific incentives to potential cancelers.
- Furthermore, without data analytics, estimating a person's likelihood of canceling a reservation would have been a vague idea.
- To determine cancellation likelihood, the logistic regression was chosen as the final prediction model.
- The performance of logistic regression improves dramatically as the number of the training data grows. Because our training dataset is so enormous, this is most likely what happened in this situation.

Based on validation misclassification rate:

- **Logistic Regression: Best model**
- **K-NN: Worst model**

# Managerial Insights

Hotel industry tend to predict cancellations based on advance booking information (aka pick-up models- classical and advanced, where they see that several bookings is picked up from a specific time point to another) for hotel demand forecast.

Classical pick-up model only utilizes completed bookings in their forecasting.

Advanced pick-up method uses both complete and incomplete bookings.

Managers can use our findings in following ways:

- Predict the possibility of cancellations; Incentives can be offered to potential cancellers to attract them to visit the hotels.

- Consider advanced pick-up method (Lead_time being an important variable in our findings)

Some suggestions include:

- Provide a 5% discount, for example, to all possible cancellers.

- On the basis data analysis, a more efficient way would be to offer targeted incentives to potential cancellers, some possible incentives can be provided, keeping in mind the availability of resources and capital when giving incentives:
- Attract travelers who book Room A - a 5% discount.
- $20 discount on BB.
- Travelers/customers from Portugal or Germany will be offered a 5% discount.