# CUSTOMER ANALYSIS OF THE HOTEL INDUSTRY

ISt687 – Final Project

Team 1
Leland Ball | Carlos Doria | Laura Pomerene
Spring 2020

# Table of Contents

# Project Background

This project is an exercise in analysis of Hyatt brand hotels customer survey response data. The objective of the analysis will be to provide actionable insights to management through the use of advanced modeling, with a focus on any one or all of the following objectives: improving service, improving revenue or increased utilization of amenities.

# Project Scope

The scope of this project covers customer survey responses from hotel visits between January and March 2014. The data contains 19,342 records located in the Americas, Europe, The Middle East & Africa, and Asia Pacific.

# Business Questions

In reviewing the data dictionary and project scope, but prior to reviewing or analyzing the data, the following business questions were established:

1.      What do our best customers look like?

2.      Are people from certain countries more or less likely to be satisfied with their accommodations

3.      Are staff at certain hotels under-performing, even when accounting for more rigorous raters?

4.      What are our top performing hotels? How would we define that? In revenue or net promoter?

5.      Which countries do our hotels do best in?

6.      What demographic is most likely to rate us higher? Are we doing a better job serving that customer?

7.      Which demographic rates us the lowest?

8.      Is there a survey metric correlated to the guest's likelihood to recommend?

9.    Does the purpose of the visit (business or leisure) affect net promoter? Do certain hotels cater to one or the other?

10.   Does price factor into a favorable net promoter rating? Does business or leisure purpose affect that?

11.   What is the profile of a promoter? Do they have other key attributes, such as travel purpose or a rewards tier?

12.   Does membership in gold/platinum/diamond rewards programs affect the overall hotel ratings?

13.   Are favorable survey responses more likely to come from those traveling domestically or from another country?

14.   Which hotels are improving or up and coming? Are any declining?

15.   Do owned or franchised hotels perform better?

## Final Business Questions

After analyzing the data set and determining our valid data points, we narrowed our business questions down to 6 to drive our business insight.

1.  What do our best customers look like?

2.  What are our top performing hotels? How would we define that? In revenue or net promoter?

3.  Is there a survey metric correlated to the guest's likelihood to recommend?

4.  Does membership in gold/platinum/diamond rewards programs affect the overall hotel ratings?

5.  Which hotels are improving or up and coming? Are any declining?

6.  Do owned or franchised hotels perform better?

# Data Quality Assessment

In our initial assessment of the data we learned that 92% of the records were for hotels located in the United States. Non-American regions did not have all 3 relationship types; managed, franchised, and owned. The timeframe of responses were from check in dates between January 2014 and March 2014. There was an outlier of a check in date of February 2013, which we removed. We ended subsetting the data to exclude hotels from other countries with stays in early 2014.

# Data Cleansing Methodology

We subsetted the data to only include hotels located in the contiguous 48 states in the U.S. In our data load function, we assigned each column their respective format; text, date, and numeric. We munged "GP_Tier" by changing their format to as.character, making all results lowercase, and replacing shorthand names with their full name. We converted the "Children_Num_C" and "F.B_Freq_H" columns to binary fields. In these cases, we interpreted 0 (zero) children as traveling without children and treated 0 (zero) F&B as to guests not having visited the food and beverage outlets. We also added a "weeks" column, based on the date. Genders that were NAs were redefined as "prefer not to answer". We did not have a column with the hotel's state. We saw that it was in the data dictionary, but not in the data set. We used latitude and longitude to create "STATE_PL". For all the survey questions, we replaced the NAs with the mean. We used the following code to remove the 36 columns we did not need:

```
# Remove unused columns
data[,c("MARKET_GROUP_C", "ROOM_TYPE_CODE_C", "WALK_IN_FLG_C", "CHECK_IN_DATE_C",
        "CHECK_OUT_DATE_C", "LENGTH_OF_STAY_C", "NUMBER_OF_ROOMS_C", "ADULT_NUM_C",
        "POV_CODE_C", "QUOTED_RATE_C", "RESERVATION_DATE_R", "ENTRY_TIME_R",
        "ENTRY_HOTEL_CODE_R", "LAST_CHANGE_DATE_R", "ROOM_TYPE_CODE_R", "STATE_R",
        "GUEST_COUNTRY_R", "PACE_CATEGORY_R", "PACE_R", "REVENUE_USD_R", "Guest_Country_H",
        "Age_Range_H", "Language_H", "Hotel.Name.Long_PL", "Award.Category_PL", "City_PL",
        "Country_PL", "Ops.Region_PL", "Currency_PL", "Dom.Int.l_PL", "Brand_PL",
        "Club.Type_PL", "Region_PL", "Category_PL", "Class_PL",
        "Booking_Channel")] <-list(NULL)
```

# Fields & Variables

The variables we selected based on logical segments that may drive or explain net promoter behavior. We selected 32 of the 58 variables to start the analysis process and answer our questions, we soon narrowed that down to 22, eliminating variables including country, region, length of stay or # of adults, when we could not determine a correlation or a significant finding that would otherwise answer a business question. 6 additional fields were added to better fit a modeling exercise, perform aggregate analysis or more easily plot data points.

| Field | Data Type | Description |
|---|---|---|
| CHILDREN_NUM_C | numeric | # of children present during stay |
| Gender_H | text | Guest's gender [male; female; prefer not to answer] |
| Likelihood_Recommend_H | numeric | Likelihood to recommend metric; value on a scale 1-10 |
| Overall_Sat_H | numeric | Overall satisfaction metric; value on a scale 1-10 |
| Guest_Room_H | numeric | Guest room satisfaction metric; value on a scale 1-10 |
| Tranquility_H | numeric | Tranquility metric; value on a scale 1-10 |
| Condition_Hotel_H | numeric | Condition of the hotel metric; value on a scale 1-10 |
| Customer_SVC_H | numeric | Quality of customer service metric; value on a scale 1-10 |
| Staff_Cared_H | numeric | Staff cared metric; value on a scale 1-10 |
| Internet_Sat_H | numeric | Internet satisfaction metric; value on a scale 1-10 |
| Check_In_H | numeric | Quality of the check in process metric; value on a scale 1 - 10 |
| F.B_FREQ_H | numeric | # of times guest visited Food & Beverage outlet |

| | | |
|---|---|---|
| **F.B_Overall_Experience_H** | numeric | Overall FB experience metric; value on a scale 1-10 |
| **Property_ID_PL** | text | Unique hotel identifier |
| **Hotel.Name.Short_PL** | text | Abbreviated hotel name |
| **Property.Latitude_PL** | numeric | Latitude coord of hotel |
| **Property.Longitude_PL** | numeric | Longitude coord of hotel |
| **Type_PL** | text | Type -Business Select Service |
| **Location_PL** | text | Location type [urban, airport, resort] |
| **Relationship_PL** | text | Relationship of the hotel with Hyatt Corporation |
| **GP_Tier** | text | GP Tier of the guest from multiple program sources [gold, diamond, platinum, other] |
| **NPS_Type** | text | Indicates the direct response from Likelihood to recommend [Promoter, Passive, Detractor] |
| **CHECK_OUT_DATE_MONTH** | text | Month of check out |
| **CHECK_OUT_DATE_DAY** | text | Date of check out |
| **WEEK** | numeric | Week number in date range |
| **STATE_PL** | text | Hotel's geographic state |
| **Hotel** | text | Hotel name |
| **imp_and_dec** | text | Improved or declined, based on NPS score. |

# Initial Analysis and Visualizations

After running our preliminary review, NPS_type, a segmentation of the Likelihood to Recommend metric, stood out as a key variable to answer our business questions. Originally developed by Bain and Company, the Net Promoter Score (NPS) is a metric derived from a single question gauging the Likelihood to Recommend. On a 10 point scale Promoters give 9s and 10s, Passives give 7s and 8s and Detractors are 6 and below. NPS is calculated by subtracting the % of detractors from the % of promoters. This simple calculation helps an organization determine where they should allocate their resources, which most often is directed at retaining existing promoters and creating new Promoters from the neutral Passives (Rathore 2019). NPS alone can often be misinterpreted, so it's important to provide context with other metrics and customer comments. According to Perceptive Research (2018), the most common complaints of hotel industry detractors are:

- Lack of value for money
- Dated decor
- Poor customer service
- Poor room hygiene

- Poor facilities
- Uncomfortable beds
- Lack of WiFi

Unfortunately, our survey data set does not include comment data, however we do have several satisfaction metrics that align with service, facility and room quality. NPS is calculated in aggregate based on a segment of the population. For this analysis we chose to apply NPS to multiple segments:
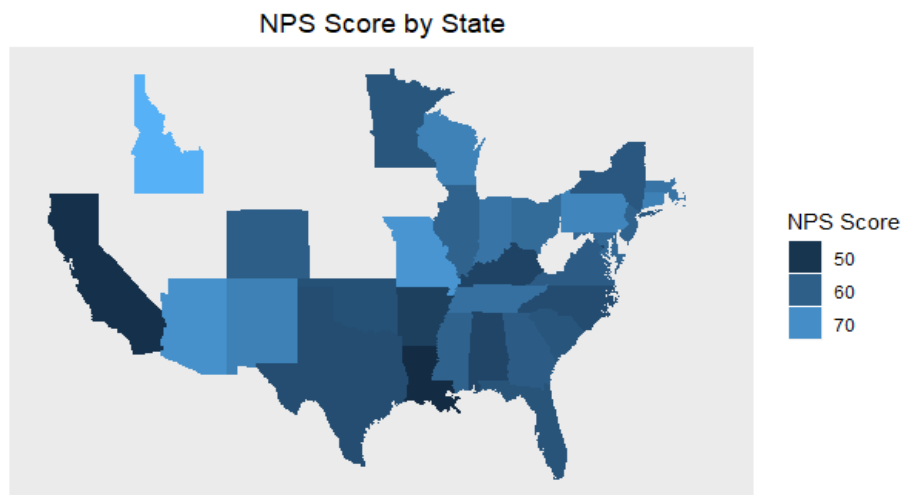
- Demographics
- Hotel Property
- Ownership Type
- Geographic area
- Change over time
- GP-Tier - rewards/affinity program level

Identification of Promoters in a population is significant in that compared to detractors, promoters are almost six times as likely to forgive when they encounter a problem, are

more than five times as likely to repurchase, and are more than twice as likely as detractors to actually recommend a company (Rathore 2019).
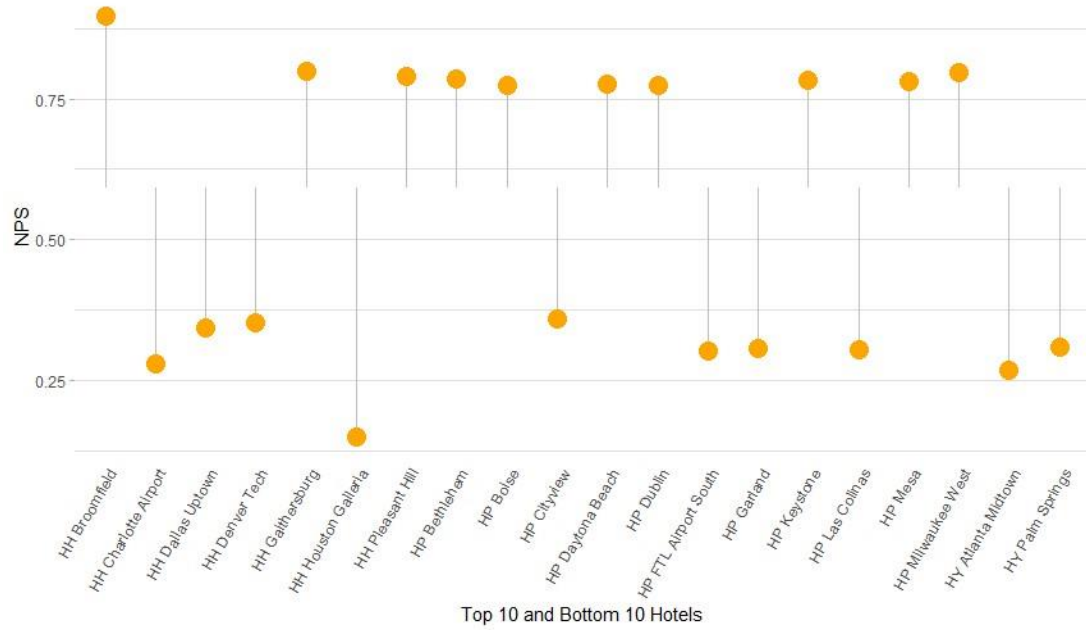
According to Bain and Co, NPS accounts for 20% to 60% of a company's organic growth rate. And, on average, the leader in an industry has an NPS more than double of its competitors (Rathore 2019). Given higher scores indicate more loyal customers, which leads to increased revenue and higher profits, we calculated NPS by segment. Some of the ways we used this are:

- **NPS by Hotel Over Time** - period 1 compared to period 2, looking for hotels that are improving or declining. This is detailed in our Statistical Score Improvement Model.
- **NPS by GP-Tier -** comparing rewards levels and the scores above or below the mean by GP-Tier to test utilizing Naive Bayes is detailed in our modeled approach to predict NPS by GP-Tier.
- **NPS by State -** visually we wanted to see if we could learn anything from a heat map of the US, indicating NPS by State. As mapped below, the states with higher scores are lighter blue.



NPS Score by State

- **NPS by Hotel -** To answer the question: Who are our top hotels? We took an average NPS score of 59 and plotted our top and bottom performers. They ranged
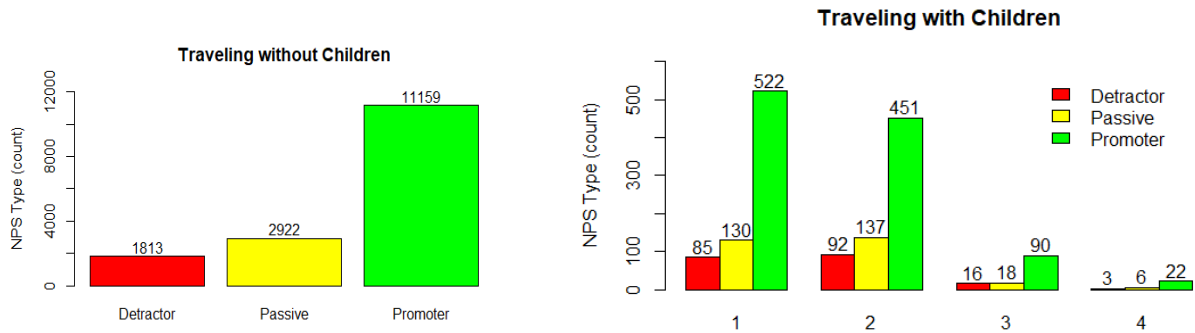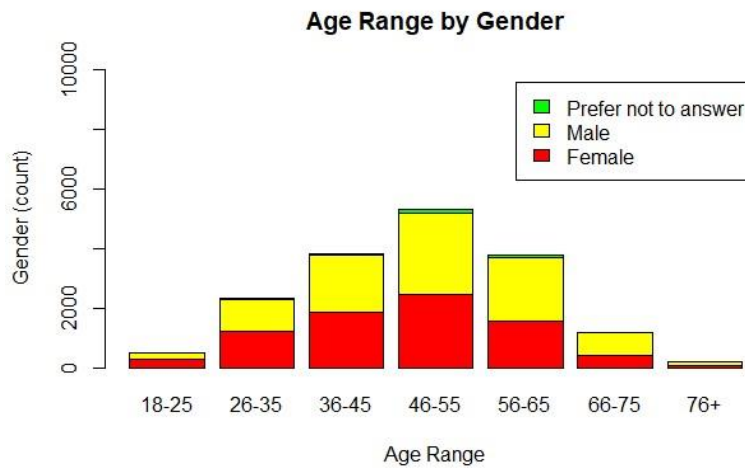
from the top hotel, the Hyatt Broomfield, Colorado at 89.7 to the lowest ranking at the Houston Galleria with a 32.5.



Top 10 and Bottom 10 Hotels

# Descriptive Statistics

The profile of the average **Hyatt Hotel customer**:

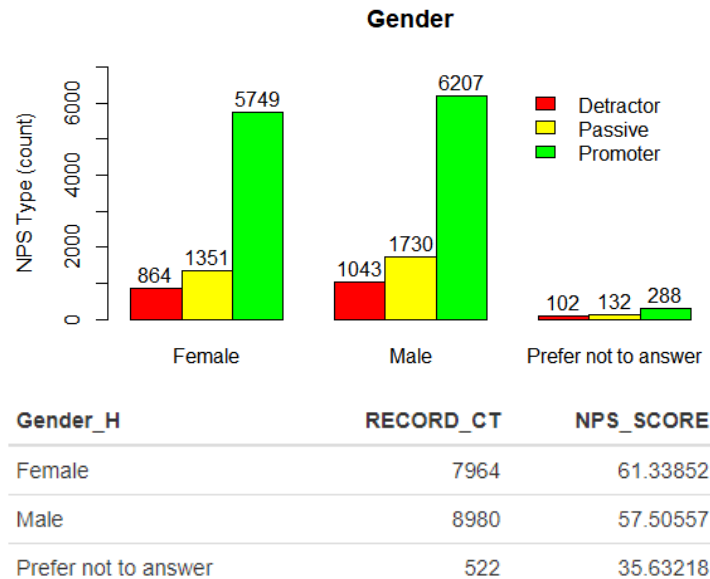- The median age range is 46-55, with 75% between ages 35-65
- Skews slightly male at 52%
    - Travels for business (70%) and most often without children



Age Range by Gender



Traveling without Children



Traveling with Children

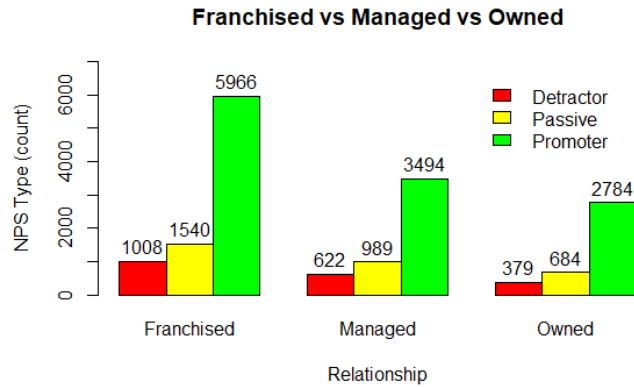| CHILDREN_NUM_C | RECORD_CT | NPS_SCORE |
|---|---|---|
| 0 | 15894 | 58.80206 |
| 1 | 737 | 59.29444 |
| 2 | 680 | 52.79412 |
| 3 | 124 | 59.67742 |
| 4 | 31 | 61.29032 |

**What does our best customer look like?**

We've defined "best customer" as one with an NPS_Type = "Promoter". The chart and table below show that female customers have a higher NPS over males by nearly 7%.



| Gender_H | RECORD_CT | NPS_SCORE |
|---|---|---|
| Female | 7964 | 61.33852 |
| Male | 8980 | 57.50557 |
| Prefer not to answer | 522 | 35.63218 |

# Interesting Findings

- Leisure travelers have a slightly higher concentration of promoters at 72% compared to business travelers at 70%.
- State_PL field is missing from the data set. This required a workaround to identify hotel locations in each state
- No amenities were found in the data (no "basket" items available for associative rules mining)
- Travelers with two children were the toughest critics of all travelers with children
- Resorts have a higher NPS than other types of hotels
- Owner vs Franchise: Owned hotels top Franchisees for higher NPS values by 7 percentage points
- Owned hotels are 11 percent higher than managed hotels when it comes to NPS
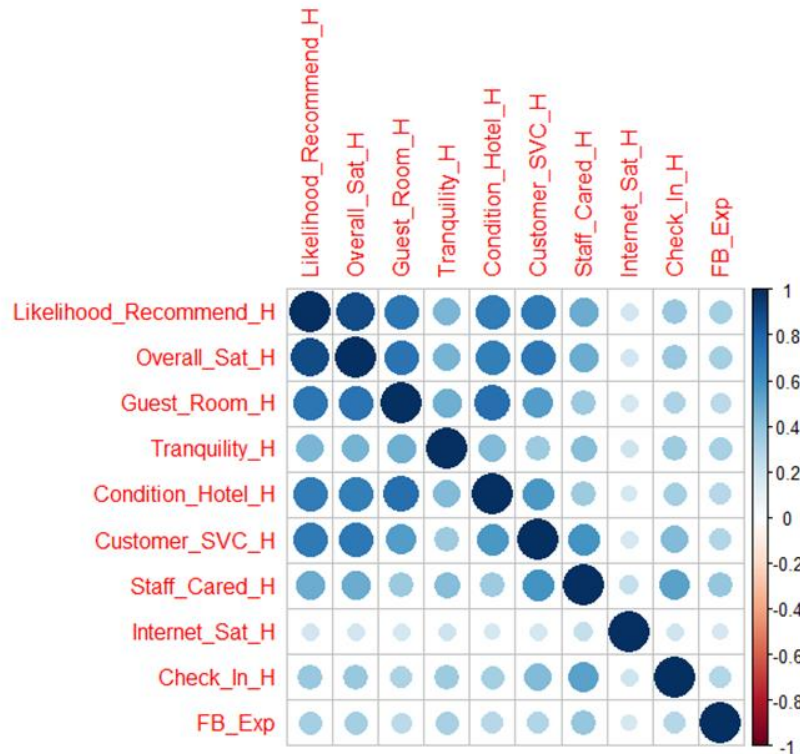
**Franchised vs Managed vs Owned**



| Relationship_PL | RECORD_CT | NPS_SCORE |
|---|---|---|
| Franchised | 8514 | 58.23350 |
| Managed | 5105 | 56.25857 |
| Owned | 3847 | 62.51625 |

# Modeling Techniques and Results

## Linear Models as Measures for Customer Priorities

Is there a survey metric correlated with a guest's likelihood to recommend a hotel? The answer to this question will affect future surveys, and inform the focus and finances for Hyatt hotels. Answering this question began with collecting the current survey metrics and running a correlation across each metric. The data for most of these fields ranges in value from 1 to 10, with missing values filled in with the mean value for that field. To begin with, a correlation matrix was first developed, to guide the further development of multiple linear models. This correlation was run over the entirety of the data set, and the resulting correlation matrix can be seen below.

In the plot above, larger circles and darker colors both imply stronger correlation coefficients. The correlation analysis shows that all survey fields are positively correlated, though some are more strongly correlated than others. Of particular interest is the Likelihood_Recommend_H field. It is this field that is used to determine if a guest is considered a promoter, a passive, or a detractor of the Hyatt hotel that they rated. The three categories are further used to calculate the NPS values for each hotel, which is used extensively in the analyses.

Next, a comprehensive linear regression model was created, which showed several survey questions to play a significant role in affecting the likelihood that a guest would recommend the hotel, as seen in the chart below.

| Survey Category | Relationship | Significance |
|---|---|---|
| Guest Room Satisfaction | directly | *** |
| Customer Service | directly | *** |
| Hotel Condition | directly | ** |
| Internet Satisfaction | directly | ** |

| Staff Caring | inverse | |
|---|---|---|
| Check In Process | directly | |
| Food and Beverage Overall Experience | directly | |
| Tranquility | directly | |

This comprehensive linear model had an adjusted R squared score of 0.82, which means that it explains 82% of the variation found within the Likelihood_Recommend_H field. While the linear model's terms for guest satisfaction with their rooms, the customer service, and the hotel condition all make intuitive sense, there are a couple anomalies worth mentioning.

The metric for Staff Caring is the only field that has an inverse relationship. This would be an unusual occurrence because it implies that as people rated their satisfaction with the amount of care the hotel staff gave them, their likelihood to recommend the hotel as a whole decreased. This is not relevant, because the significance for the relationship of this field is lower than the 95% or even 90% confidence interval that this model was made with.

Another interesting relationship is the significance that Internet Satisfaction played in this model. While the correlation matrix showed this field to be one of the least significant, the more involved regression analysis says otherwise. To get a more in-depth look at how well this relationship explains the data, a series of individual linear models were made, to take a look at each survey field. The results are charted below.

| Survey Category | R-Squared Score for Single-Term Model |
|---|---|
| Guest Room Satisfaction | 0.672 |
| Customer Service | 0.589 |
| Hotel Condition | 0.566 |
| Internet Satisfaction | **0.136** |

While the individual models show a large degree of explanatory power for the satisfaction of guest rooms, customer service, and hotel condition, the same cannot be

said for the Internet model. While this variable may have some explanatory power, only about 1% of the model's explanatory power was lost upon removing this term.

After removing internet satisfaction from the model, the adjusted R-Squared value became **0.806**. Furthermore, the other measurements of model fitness, the F-statistic and the p-value associated with the model as a whole indicate that these three terms alone can successfully be used to predict a large amount of the variability behind a customer willing to recommend the hotel to others.

The purpose behind this modeling effort is less about developing a model to predict the answer to a customer's survey question from some other survey questions, as it is to find what areas truly drive overall satisfaction in a guest's stay. From the models developed from the data, it is clear that a guest's primary concerns are for the condition of the hotel, satisfaction with their rooms, and customer service. These three areas should be prioritized for the purposes of allotting time, training, and finances.
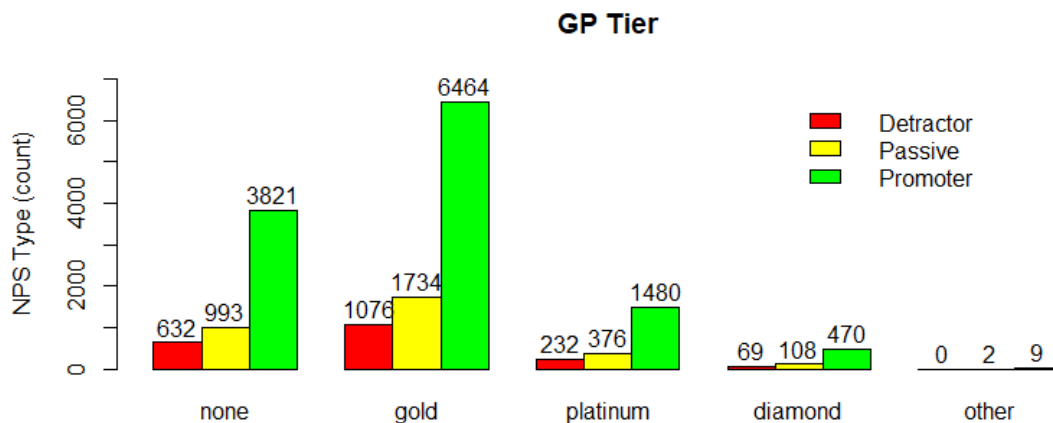
The merits of these survey questions can also be seen in this analysis. If the primary metric that the hotel chain prioritizes is the likelihood of a customer to recommend their hotel to others, then many of the other survey questions are found to be ineffective predictors of this most important metric. If a new survey is to be developed, this research should be taken into account so that more relevant questions can be asked of guests in the future.

## Predicting NPS by GP Tier With Multiple Models

Does membership in a GP Tier increase the likelihood of a customer promoting a hotel? The GP Tier system consists of three main customer loyalty categories that receive different rewards and benefits. These categories include gold, platinum, and diamond. The answer to this question will begin to inform us on whether or not the GP Tier system can be used as a tool to increase customer promotion of the hotel chain.

To begin, the tiers were explored by taking the NPS values calculated from all customers in each tier, in addition to those customers not in a tier. As mentioned earlier, the NPS values were derived from the Likelihood_Recommend_H survey field.

| GP Tier | Record ct | Mean NPS |
|---------|-----------|----------|
| diamond | 647 | 0.62 |
| platinum | 2088 | 0.60 |
| gold | 9274 | 0.58 |
| none | 5446 | 0.59 |



The results shown above do not hold great hope for discovering a relationship between GP tier and a customer's likelihood to promote the hotel, as the promoter scores are not very different from one another. Venturing further, a linear model was next developed, to better understand if there was such a relationship.

In order to run this categorical data through a linear model, it must first be dummy coded into many fields, each coding for the binary presence of a tier. Customers without a tier are encoded as zeroes in all other tier columns. The model was then calculated and returned an Adjusted R-Squared Score of 0.449. Not only is this score somewhat low, implying that only about fifty percent of the variation is explained by this model, but the other measures of model applicability (the F-Statistic and p-value of the model) describe this model to be a poor predictor of a favorable hotel rating, and must be discarded.

Not to be dissuaded, a last attempt at modeling the behavior between GP Tier and a customer's likelihood to recommend was attempted in the form of a Naive Bayes classifier. The prerequisites for this model includes using the original categorical version of the recommendation likelihood column, with the detractors, passives, and promoters

categories. The model was trained on two thirds of the data, and tested for fitness using the remaining third.

| | **Predicted** | | |
|---|---|---|---|
| **Correct: 84.3%** | **Detractor** | **Passive** | **Promoter** |
| **Detractor** | 0 | 0 | 689 |
| **Passive** | 0 | 0 | 1070 |
| **Promoter** | 0 | 0 | 4063 |

(**Actual** labels the rows Detractor, Passive, Promoter)

The confusion matrix shows that while the Naive Bayes model is able to predict correctly 84% of the time, the model is only ever predicting that a customer will be a promoter, since the training data contained mostly promoters. Successive attempts to train the model using a balanced dataset with equal numbers from each NPS Type met with accuracies in the low double-digits. The conclusion drawn in this case is that there is no clear relationship between GP Tier and whether or not a guest will recommend the hotel to someone else.

This result is not altogether unsatisfactory, because while no model has been created to correctly predict how being in a tier will affect a review, this information can be used to guide hotels away from certain business practices. These results show that simply elevating a guest into a special tier will not garner a higher rating on survey data. Additionally, it may imply that tier amenities are not currently enticing enough for detractors to look away from the faults that they find when rating a visit.

## Statistical Score Improvement Model

Which hotels have improved over the course of this timeframe? Each hotel's performance (as rated by their guests) can be quantified in aggregate using a net promoter score (NPS). For this study, the mean NPS value for each hotel was calculated on a per-day basis. In order to determine which hotels showed a score improvement during this time, several hurdles had to be overcome.
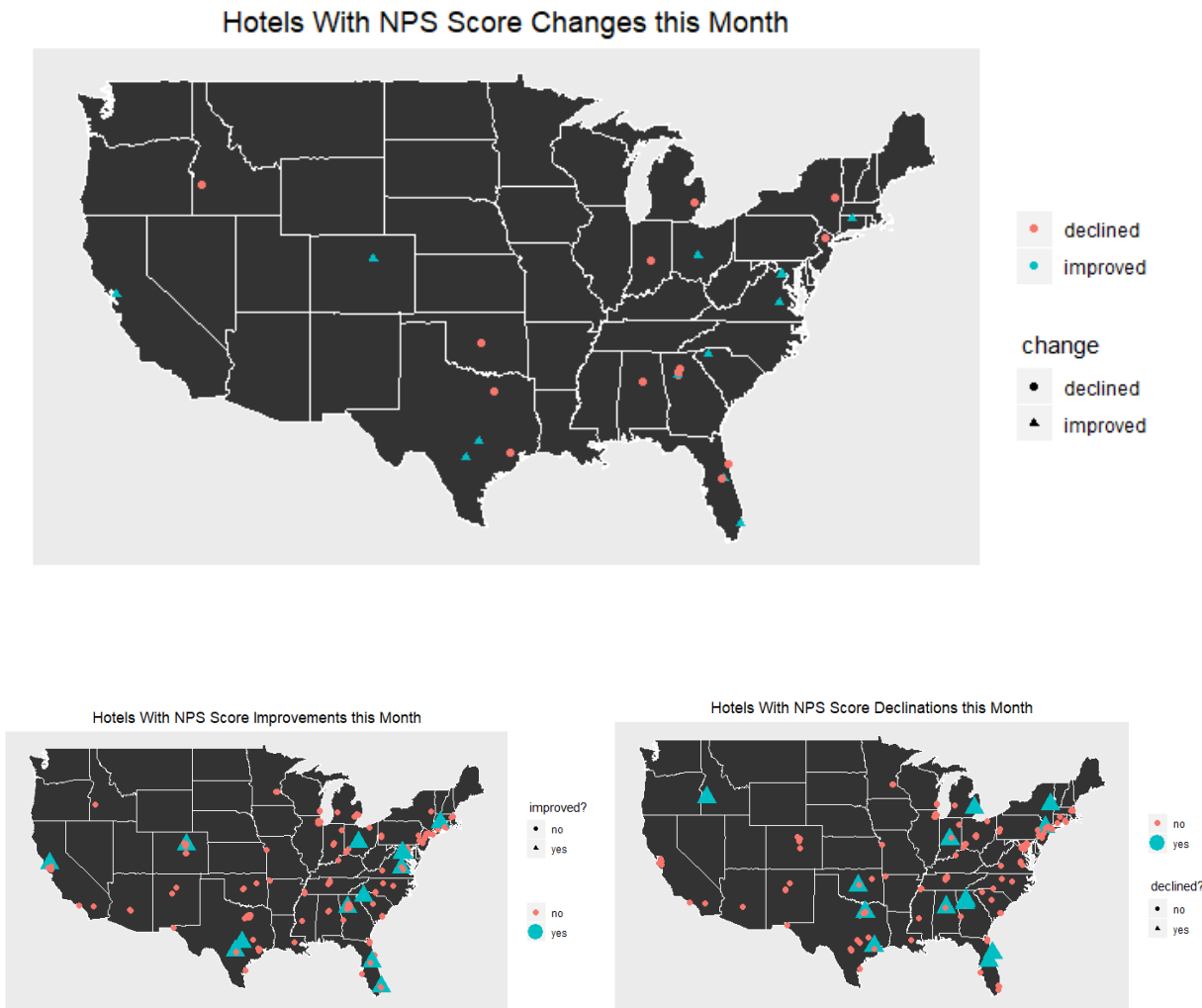
The data showed some gaps in time. Many hotels had three weeks worth of surveys, but some had four, five, or even only one week's worth of data. In an ideal world, there would be data throughout the time period inspected, for every hotel. To overcome this obstacle each hotel's data was ordered chronologically by day, and then cut in half. This resulted in two lists of scores that could be compared.

The comparison itself posed a challenge, since the data could not simply be compared to the mean NPS value for each hotel, and then subtracted. Indeed, every hotel would show some amount of variation and there would be no way of telling if the change was significant. Measuring the significance of a hotel's change was also of prime concern.

This challenge was overcome by applying a one-tailed, two-sample t-test to each hotel's "before" and "after" data. This test was performed twice: once to search for hotels showing score improvements, and once to show hotels that were declining in score. This output is included in the table below.

| Improved Hotels | Declining Hotels |
|---|---|
| HH Pleasant Hill | HP Daytona Beach |
| HP Denver Airport | HY Atlanta Midtown |
| HP FTL 17th Street | HH Whippany |
| HP Atlanta Downtown | HP Inverness |
| HP San Antonio NW | HP Boise |
| HP Sterling | HP Indianapolis Airport |
| HP Haywood | HP OKC Airport |
| HP Richmond Arboretum | HP Windward |
| HP Orlando Airport | HP Orlando Convention |
| HH Austin Arboretum | HH Houston West |
| HH Windsor | HH Richardson |
| HP Herndon | HP Auburn Hills |
| HP Columbus OSU | HP Atlanta Perimeter |

To gain a further understanding of whether geography plays a role in the NPS changes seen in the chart, various maps were constructed to show those hotels that have improved, and those in decline



Hotels With NPS Score Changes this Month



Hotels With NPS Score Improvements this Month



Hotels With NPS Score Declinations this Month

As can be seen from the maps, there are no easy conclusions that can be drawn from plotting each hotel's location. While many hotels with score improvements are in large urban areas, many of the hotels in decline are from similarly large urban areas. Many states have both improving and declining hotel ratings, which again makes it difficult to compare.

There is a path forward. Using this data, the hotel chain can further analyze business practices of the successful hotels while directing resources to the hotels in decline, if those hotels are not otherwise in need of consolidating. Additionally, this process is a repeatable metric that can be used to test the improvement of any number of hotels, with any size data. Whether they have a few weeks of surveys, or an entire quarter's worth of data to analyze. This metric can be repeatedly applied in an ongoing fashion for continuous monitoring of hotel fitness.

## Actionable Insights

The following recommendations have been developed based on our general analysis, model conclusions, inferences of national research and with a prioritization on maximizing resources for the greatest improvement in revenue.

- Hotels with declining NPS should be an area of focus, particularly in areas that most drive customer satisfaction
- Three areas that should receive continual focus are the ones that have shown a strong correlation to drive NPS: Customer service, Guest Room Satisfaction, and Hotel Condition
- Membership in GP Tier is an ineffective "carrot" to create or grow the promoter segment, and should not be used as such. Resources spent on this loyalty/rewards program should be re-evaluated for improvement to align with concerns that drive customer satisfaction, or scaled back
- Use top-performing hotels to garner best-practices to develop effective service and guest room satisfaction programs to improve underperforming hotels
- Revisit the customer experience survey questions. Consider tracking usage of amenities to predict purchase behavior with associative rules mining, and possibly tracking comments to add context to the overall metrics.

# References

Rathore, Nirmal. "The Insider's Guide to Net Promoter Score in the Hospitality Industry." *Https://Xperium.ai/2019/08/27/Net-Promoter-Score-in-Hospitality-Industry/*, Perium by Repup, 27 Aug. 2019, xperium.ai/2019/08/27/net-promoter-score-in-hospitality-industry/.

"What Is a Good Hotel Score for the Hotel Industry?" *Perceptive*, Perceptive Group, LTD, 15 Mar. 2018, www.customermonitor.com/blog/what-is-a-good-nps-score-for-the-hotel-industry.

# Appendix – R Code

# see included fine: data-analysis.R