# Group 12's Investor Report for AirBnB's in San Diego



Chinmay Gupta, Shruti Sharma, Siddhita Bhagwe, Harsh Sharma, Mansi Kosamkar, Anuj Doshi

# **Table of Contents**



# Maggle Overview

# 1. Cleaning

Handling of missing values. Numeric values were replaced with median of the particular column.

### 3. Modeling

We started with basic linear and logistic models and slowly moved onto ensemble methods.

### 2. Feature Selection

We applied domain knowledge, used feature selection models and visualizations to understand the significant variables and dropped 37 redundant variables.

### 4. Conclusion

Our final model is the XGBoost model with 29 variables which gave us an accuracy of 84% and roc estimate of 0.91

.metric	.estimator	.estimate
<chr></chr>	<chr></chr>	<dbl></dbl>
roc_auc	binary	0.913636

1 row

Confusion Matrix and Statistics

high\_booking\_rate predictedClass 0 1 0 19517 2483 1 2350 5968

Accuracy: 0.8406

95% CI : (0.8364, 0.8447)

No Information Rate : 0.7213 P-Value [Acc > NIR] : <2e-16

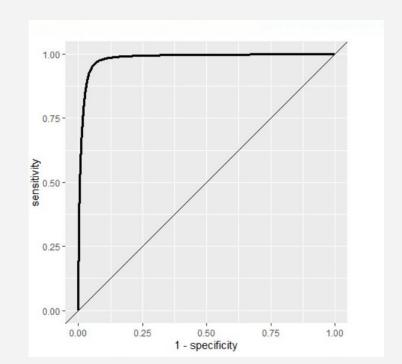
Kappa : 0.6016

Mcnemar's Test P-Value : 0.0576

Sensitivity : 0.7062 Specificity : 0.8925 Pos Pred Value : 0.7175 Neg Pred Value : 0.8871 Prevalence : 0.2787 Detection Rate : 0.1968

Detection Prevalence : 0.2744 Balanced Accuracy : 0.7994

'Positive' Class : 1



Business
Requirements

Develop a business case for an investor who is interested in acquiring homes in San Diego to put them up as AirBnB rentals.



### **Requirement 1**

Investor aims to purchase those properties which yield high airbnb traffic.



# Requirement 3

Investor needs to provide the amenities which result high airbnb traffic.

## **Requirement 2**

Investor needs to strategize how to advertise the property for airbnb rentals based on different variables.







## **Tourism**



- SAN DIEGO COMIC-CON
- SAN DIEGO COUNTY FAIR
- MIRAMAR AIR SHOW
- KAABOO DEL MAR



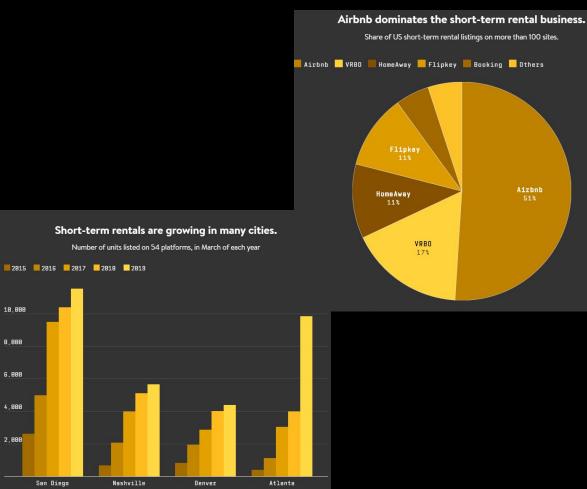
### **Popular Attractions**

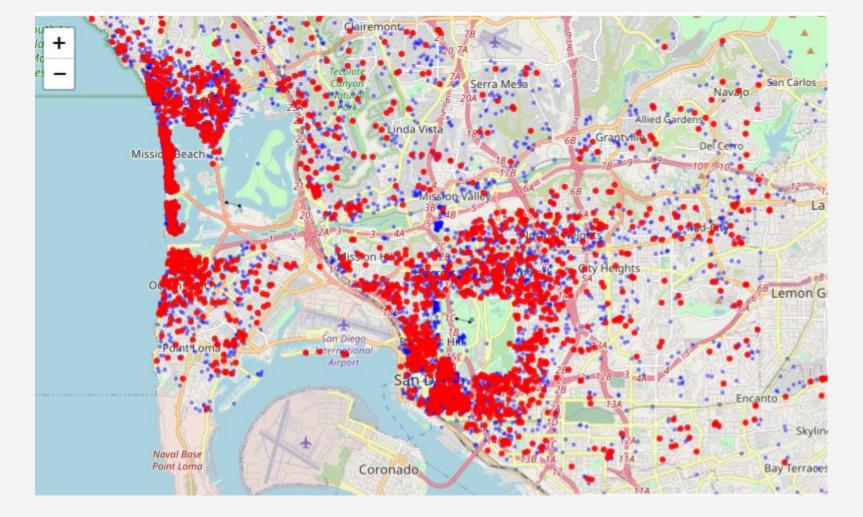
- LEGOLAND
- San Diego Zoo
- Balboa Park
- Sea World



# **Market Data Analysis**

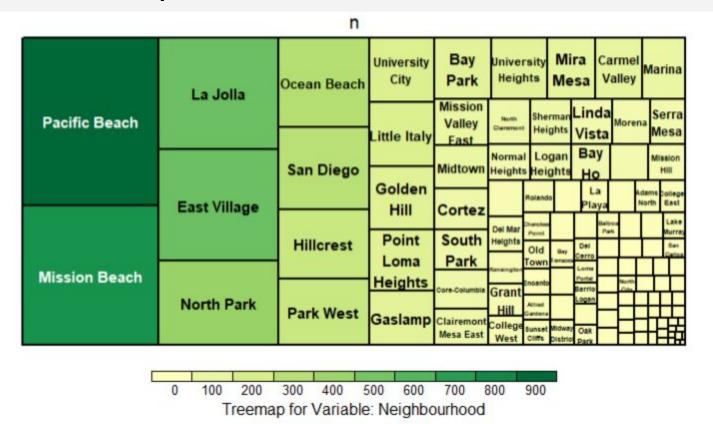






### Properties with high\_booking\_rate=1 property\_type Rancho Santa Fe Aparthotel Dome house ana Beach Apartment Farm stay Bed and breakfast Guest suite Boat Guesthouse Boutique hotel Hostel Bungalow House Loft Cabin Camper/RV Other Casa particular (Cuba) Serviced apartment Cave Tiny house Chalet Townhouse Bonita Chuia sta Condominium Villa Cottage

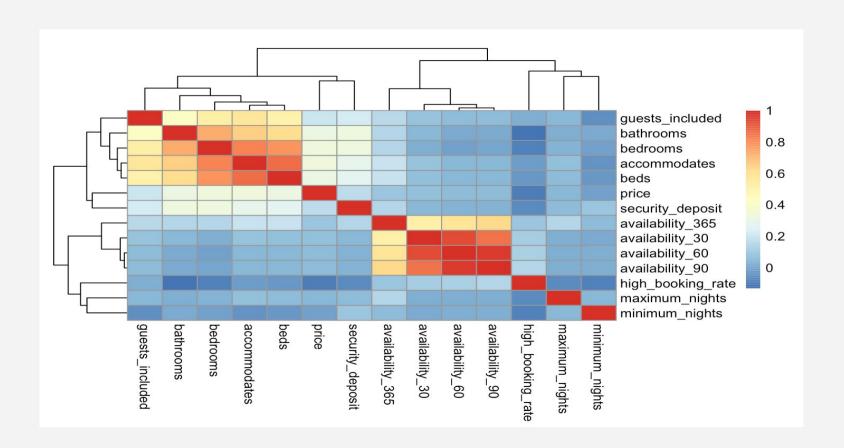
### **Tree Map**



# 1 Business Cases



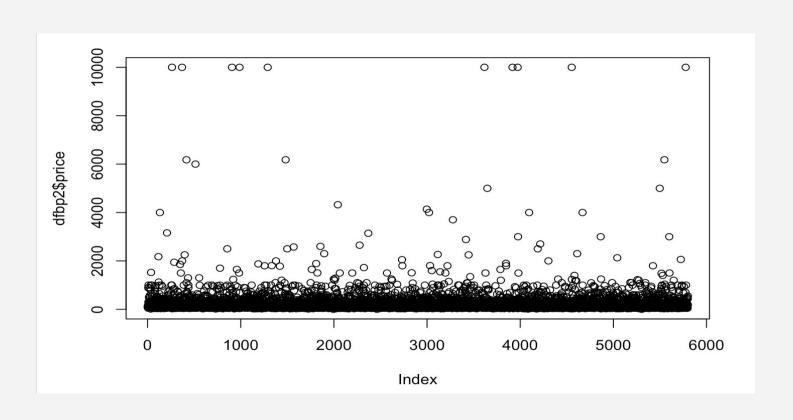
# To Focus on Marketing and Advertising an AirBnB property - We establish correlation among variables



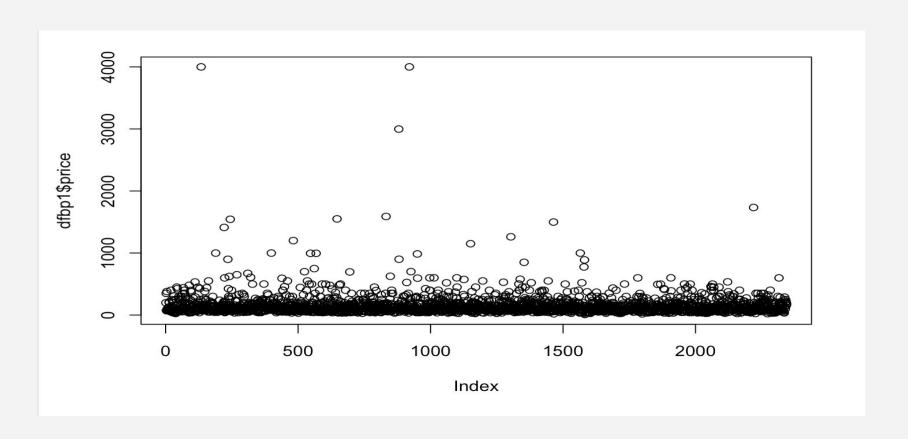
There does not exist positive correlation between Prices and High Booking Rate. So no direct inferences can be made

	high_booking_rate	price
high_booking_rate	1.0000000	-0.1011064
price	-0.1011064	1.0000000

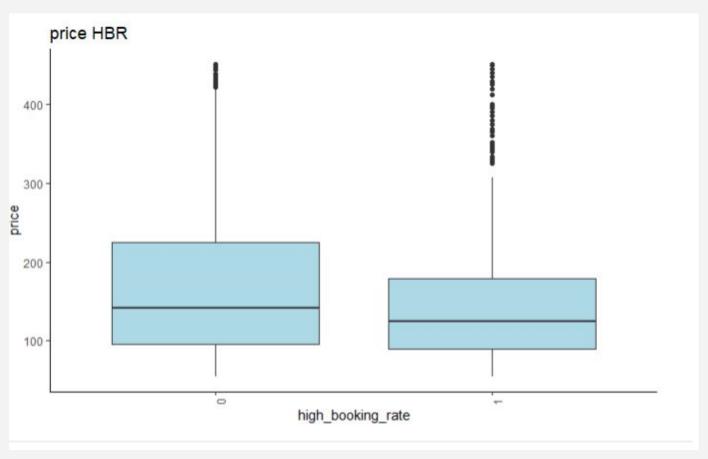
# Prices of the property with low booking rate



# Prices of the property with High booking rate



# Boxplot of prices with High booking rate



# **Conclusion:**

The best price range of the properties with higher

booking rates is under \$500.



# Variables considered

High\_booking\_rate,accommodates,availability\_30, availability\_365,availability\_60,availability\_90,bathr ooms,bedrooms,maximum\_nights,minimum\_nights, guests\_included,beds

# **Model Applied**

logistic Regression Model

```
Call:
glm(formula = high_booking_rate ~ ., family = "binomial", data = dfcTrain)
Deviance Residuals:
   Min
                 Median
                                     Max
-1.4886 -0.8532 -0.6254 1.1839
                                  3.0063
Coefficients:
                  Estimate Std. Error z value Pr(>|z|)
              -4.521e-01 9.449e-02 -4.784 1.71e-06 ***
(Intercept)
accommodates
             7.379e-02 2.518e-02 2.931 0.003378 **
availability_30
                3.212e-02 9.123e-03 3.521 0.000431 ***
availability_365 7.256e-04 3.320e-04 2.186 0.028830 *
availability_60 -6.723e-02 9.400e-03 -7.153 8.52e-13 ***
availability_90 4.182e-02 4.603e-03 9.085 < 2e-16 ***
bathrooms
          -6.147e-01 7.434e-02 -8.268 < 2e-16 ***
         -1.700e-01 5.764e-02 -2.949 0.003189 **
bedrooms
maximum_nights -3.586e-04 6.202e-05 -5.781 7.40e-09 ***
minimum_nights -7.546e-02 9.032e-03 -8.354 < 2e-16 ***
quests_included 5.316e-02 1.801e-02 2.951 0.003168 **
                3.632e-02 3.711e-02 0.979 0.327802
beds
Signif. codes: 0 '*** 0.001 '** 0.01 '* 0.05 '.' 0.1 ' 1
(Dispersion parameter for binomial family taken to be 1)
   Null deviance: 6288.9 on 5293 degrees of freedom
Residual deviance: 5736.9 on 5282 degrees of freedom
AIC: 5760.9
```

Number of Fisher Scoring iterations: 6

p-values indicates that statistically important variables are

accommodates, availability\_30, availability\_365, availability\_60, availability\_90, bathrooms, bedrooms, maximum\_nights, minimum\_nights, guests\_included

# Accuracy of the model

Confusion Matrix and Statistics

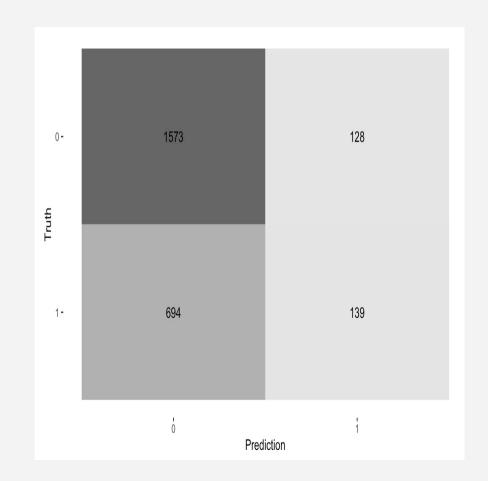
high\_booking\_rate
predictedClass 0 1
0 1573 694
1 128 139

Accuracy: 0.6756

95% CI: (0.657, 0.6938)

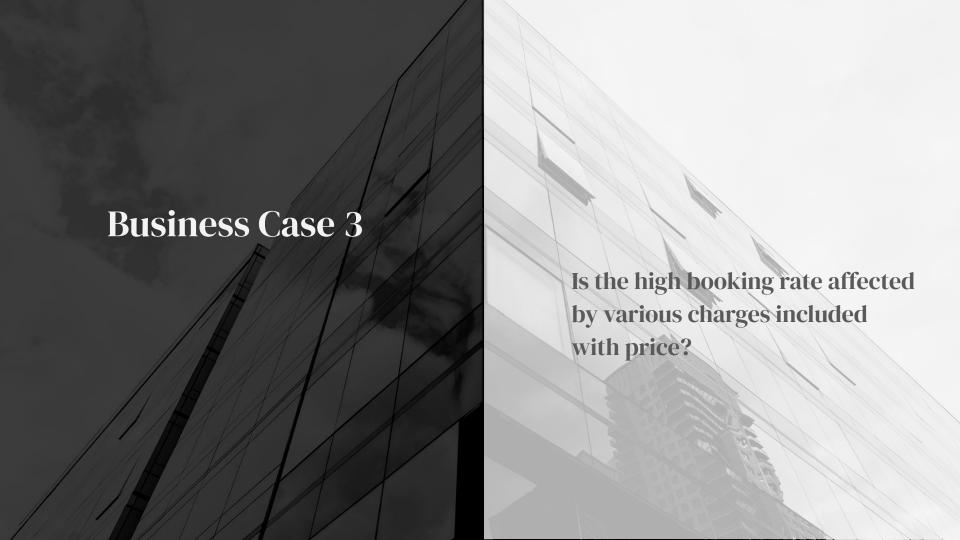
No Information Rate: 0.6713

P-Value [Acc > NIR] : 0.3292



Confusion matrix is used to find the Accuracy as the dataset was highly skewed.

**Conclusion:** For an Airbnb property to have higher booking rates, it must be flexible with respect to its booking duration and be able to provide accommodation to the guests.



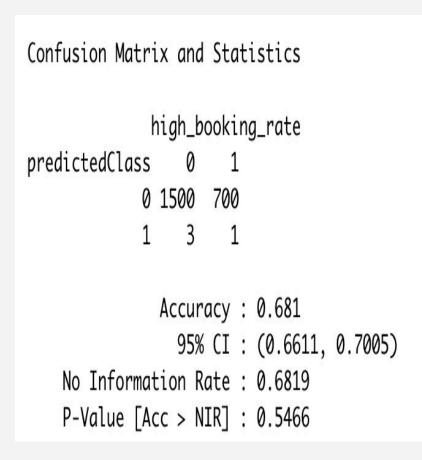
# Variables considered

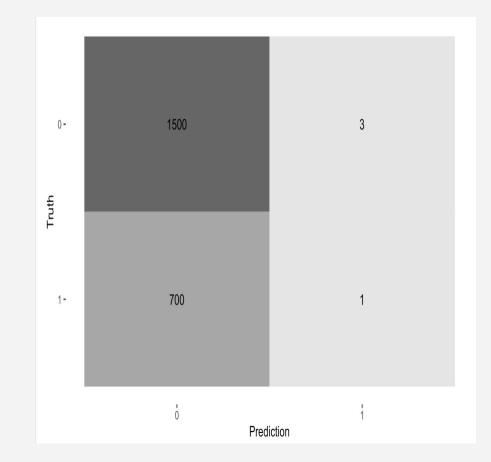
- Cleaning\_fee
- Extra\_people
- Price
- Security\_deposit

# **Model Applied**

Logistic Regression Model

```
Call:
glm(formula = high_booking_rate ~ ., family = "binomial", data = dfcTrain)
Deviance Residuals:
                                                                        Only extra-people
   Min
            10 Median
                                    Max
                             30
                                                                        and price are
-1.4240 -0.8881 -0.7884 1.4429 3.7715
                                                                        statistically important
Coefficients:
                                                                        variables
                 Estimate Std. Error z value Pr(>|z|)
(Intercept) -5.492e-01 4.877e-02 -11.261 < 2e-16 ***
cleaning_fee -8.307e-04 4.865e-04 -1.707 0.087733 .
extra_people 4.464e-03 1.190e-03 3.752 0.000176 ***
price -1.630e-03 2.554e-04 -6.384 1.72e-10 ***
security_deposit -5.069e-05 8.361e-05 -0.606 0.544305
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' '1
(Dispersion parameter for binomial family taken to be 1)
   Null deviance: 6398.2 on 5293 degrees of freedom
Residual deviance: 6240.1 on 5289 degrees of freedom
AIC: 6250.1
Number of Fisher Scoring iterations: 6
```





### Conclusion

- If a property provides for extra\_people despite charging for the same its booking rate improves.
- Cleaning fee and security deposits do not have any impact on the higher booking rates.
- If the price of a property is high, its booking rate decreases.



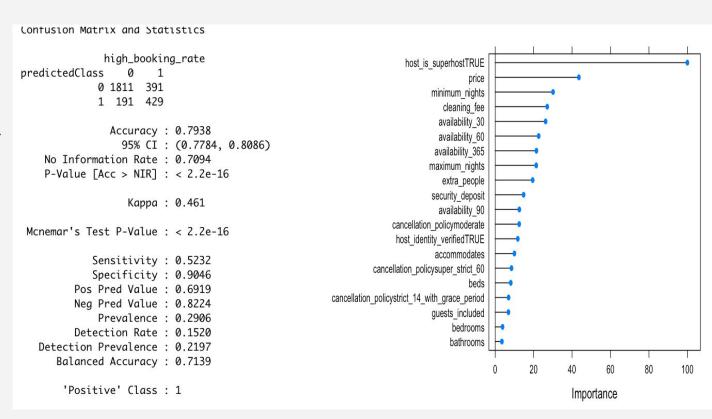
### **Variable Selection**

- Factors that contribute towards the location setting, generalized pricing strategy and in house services provided, were considered while structuring this business case
- Initially, a logistic model was run to understand the importance and significance of each variable
- Accuracy achieved for this model was 77.49%

```
Deviance Residuals:
-2.2849 -0.7233
Coefficients:
                                                 Estimate Std. Error z value Pr(>|z|)
(Intercept)
host_is_superhostTRUE
price
minimum_nights
cleaning_fee
availability 365
availability_60
availability_30
maximum_nights
extra_people
security deposit
availability_90
cancellation_policymoderate
host identity verifiedTRUE
                                                5.520e-01 6.833e-02
accommodates
cancellation_policysuper_strict_60
cancellation_policystrict_14_with_grace_period
quests included
bedrooms
bathrooms
                                               -3.396e-01 7.541e-02 -4.504 6.68e-06 ***
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
(Dispersion parameter for binomial family taken to be 1)
    Null deviance: 6870.6 on 5700 degrees of freedom
Residual deviance: 5340.6 on 5680 degrees of freedom
AIC: 5382.6
Number of Fisher Scoring iterations: 6
```

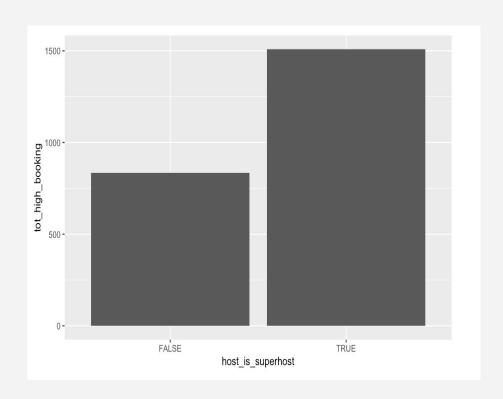
### XG Boost as a reference model for variable selection

The Xgboost model showed an accuracy 79.38% and Sensitivity of 52.32%. The sensitivity is more this important particular case as we don't want to overestimate the High\_booking\_rate factor.



# How does a host being a superhost affect booking rate?

The most important variable from the XGBoost model is the host\_is\_superhost = True. From the bar graph, we can see that the booking rate for host\_is\_superhost = True is about 55% more than where host is not a superhost

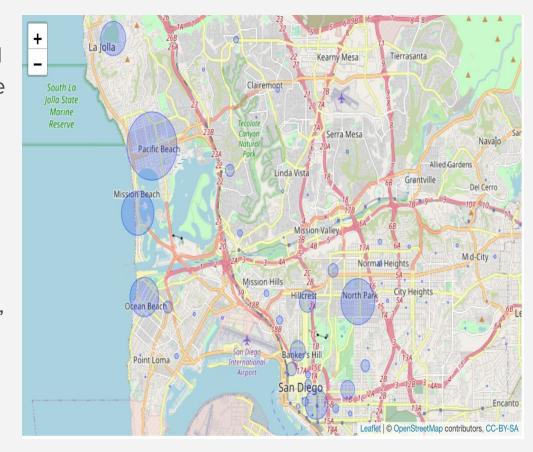


### Conclusion

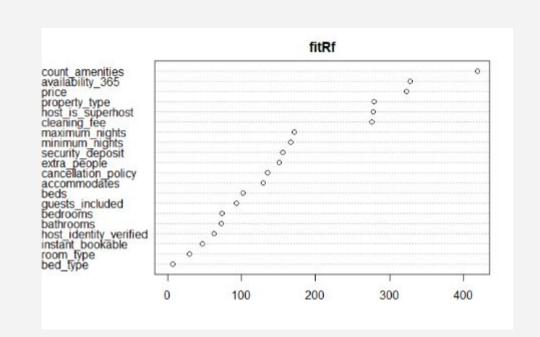
- Variables, excluding security deposits, beds and bedrooms are statistically important
- While the host holding a superhost status improves the booking rate the most

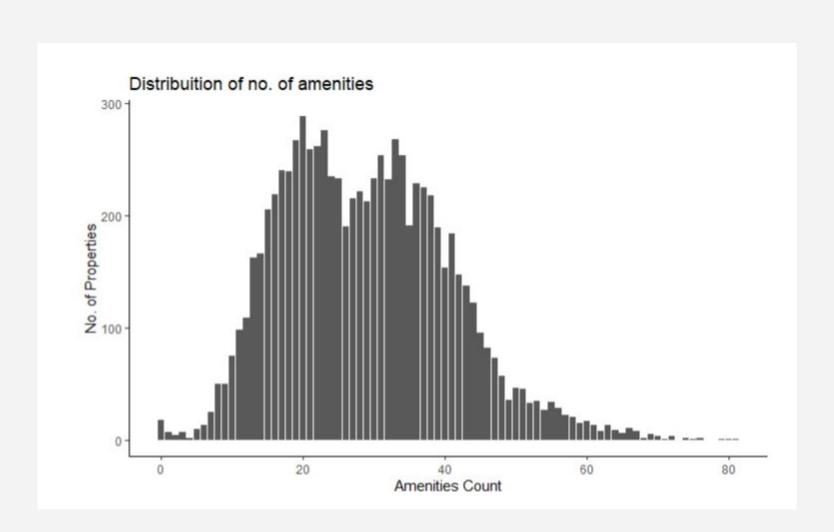


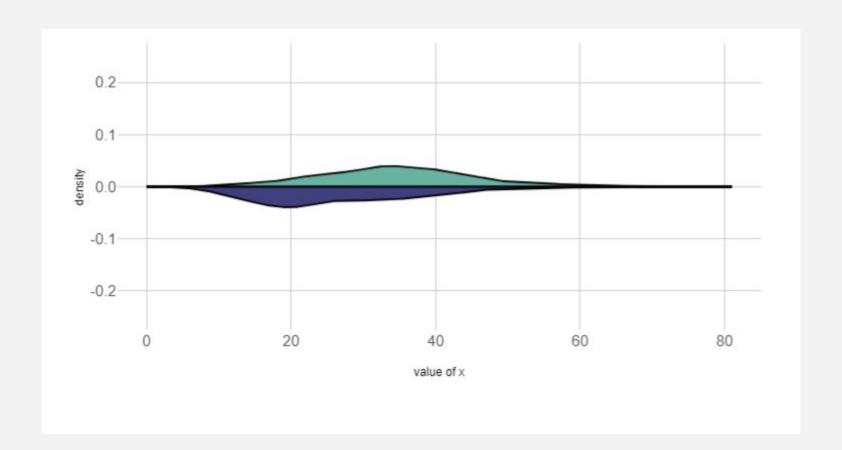
- The regions and the number of high booking rates are displayed on the map. The size of the circle also correlates with the number of high booking rates. These regions can be used as a factor to decide on initial acquisition as well as pricing.
- Due to various events like the Comic Con, Surfing tournaments, and La Jolla festival at La Jolla and general attraction for beaches, we assume that these regions have a higher booking rate.











extinguisher microwave blinens heating stove first iron allowed fire kid water hot dryer hair ai basics cable private detector self in maker laptop parking wifi coffee lock lock balcony of friendly shampoo bed aid steps long workspace



