

Marketing Analytics: Airbnb Final Project - Hawaii

Group 6: Aditi, Airi, Pankhuri, Pooja

Quick Recap: Hawaii, United States

In our project, we picked **Maui**, which includes the most popular and expensive neighborhoods, Kihei and Lahaina. Maui also has the highest average price per night in Hawaii.

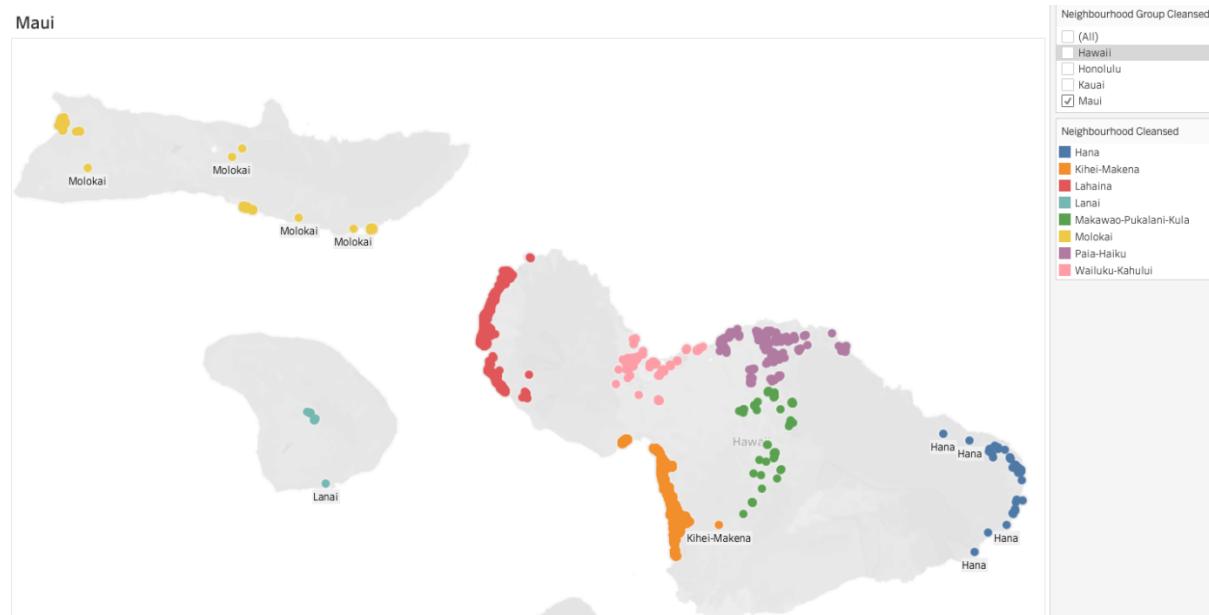
As a **guest** of Airbnb, they would like to know the **average price per night**, which is the phenomenon of interest (DV/y variable). Guests may have questions, like how much is the average price per night in a certain neighborhood with certain number of people, room type, and the number of bedrooms. These questions are important in determining and predicting the budget for the stay in Maui.

Data: We used the detailed listing data from Hawaii, Inside Airbnb

- Listings detailed
- Reviews

Based on intuition, the key variables/factors that affect and can predict your phenomenon of interest, (independent variables = x variables) are:

- *Room Type*
- *Instant Bookable*
- *Neighborhood in Maui*
- *Number of Bedrooms*
- *Number of People (Accommodates)*



Recap: Reccomendation

We concluded that all the variables influence the price, but the neighborhood influences the price the most in our model.

Kihei and Lahaina are the most visited neighborhoods in Maui and that's where we can find the most Airbnb listings. As a guest, you must make adjustments to the variables in considering the budget.

By using the profiler, we were able to predict the average price per night using these variables.

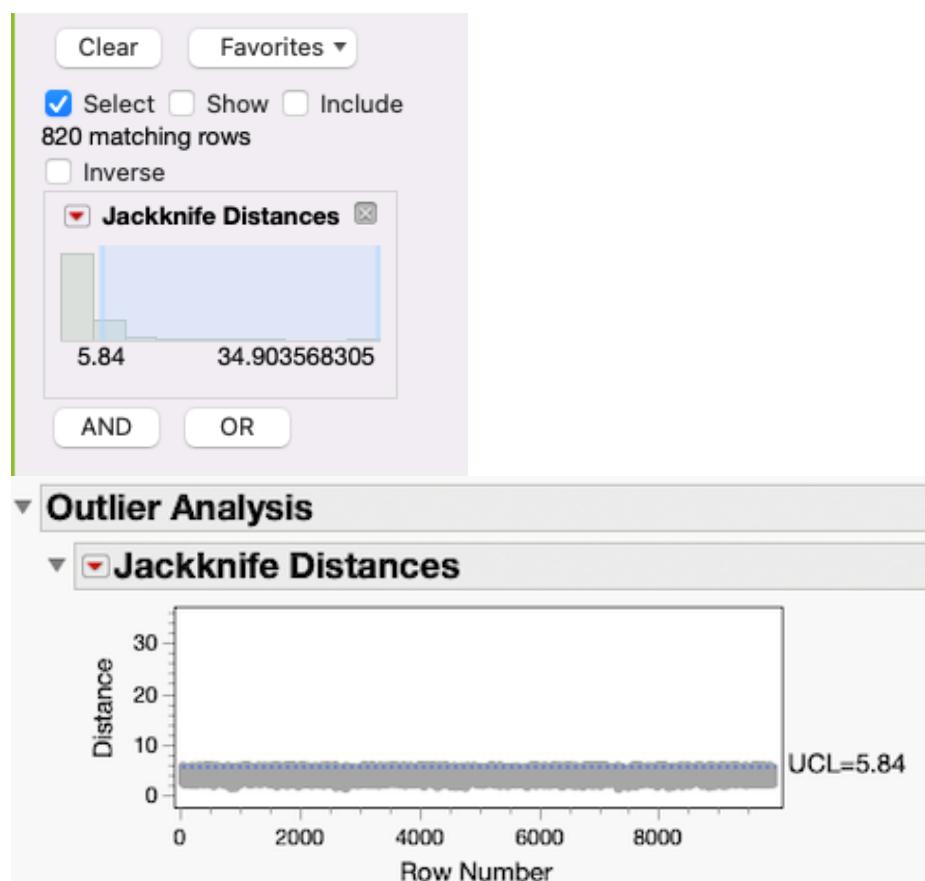
Students Vacation Trip: Private room, Instantly Bookable True, 4 people, Molokai, 2 bedroom - on average, \$270 price per night.

Couple Vacation: Private room, Instantly Bookable True, 2 people, Lahaina, 1 bedroom - on average, \$493 price per night.

Family Get-Together: Entire Home, Instantly Bookable True, 20 people, Lahaina, 10 bedroom - on average, \$18,033 price per night.

Jackknife analysis (to remove outliers):

We used jackknife analysis to remove outliers, and we excluded 820 observations.



Unsupervised Learning:

PCA/Factor Analysis:

a) Excluding your Y-variable, run the analyses with your potential continuous x-variables (the more the better)

We picked 22 potential continuous x-variables

*Number of Bedrooms
 Number of People (Accommodates)
 Bathrooms_text (number of bathroom)
 Number of Beds (bed)
 Minimum Nights
 Maximum Nights
 Number of reviews per month
 Number of Reviews
 Number of reviews Itm
 Number of reviews last 30 days
 availability_30
 availability_60
 availability_90
 availability_365
 calculated_host_listings_count
 Review_scores_rating
 Review_scores_accuracy
 Review_scores_cleanliness
 Review_scores_checkin
 Review_scores_communication
 Review_scores_location
 Review_scores_value*

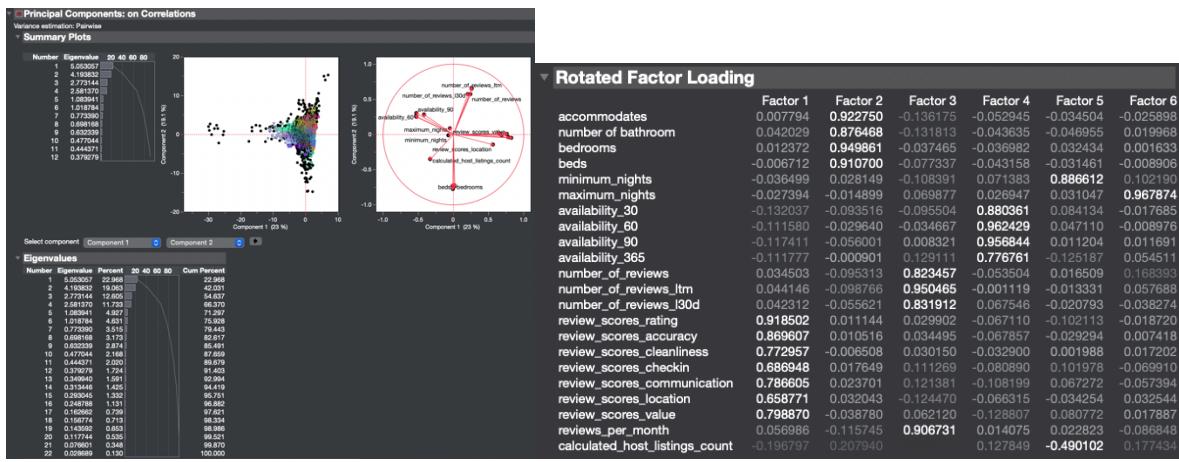
Some of these variables are highly correlated (# of bedrooms, # of bathrooms, availability, and # of reviews).

These may result in multicollinearity - variables have a higher correlation with one another and may undermine the statistically significant of an independent variable.

Multivariate												
Correlations												
	accommodates	number of bathroom	bedrooms	beds	minimum_nights	maximum_nights	availability_30	availability_60	availability_90	availability_365	number_of_reviews	number_of_reviews_itm
accommodates	1.0000	0.7688	0.8404	0.8397	-0.0140	-0.0386	-0.1092	-0.0689	-0.1040	-0.0749	-0.2031	
number of bathroom	0.7688	1.0000	0.8195	0.6986	-0.0097	-0.0149	-0.1176	-0.0704	-0.0929	-0.0515	-0.1713	
bedrooms	0.8404	0.8195	1.0000	0.8388	0.0317	-0.0081	-0.1127	-0.0609	-0.0912	-0.0386	-0.1196	
beds	0.8397	0.6986	0.8388	1.0000	-0.0161	-0.0174	-0.1100	-0.0602	-0.0906	-0.0517	-0.1590	
minimum_nights	-0.0140	-0.0097	0.0317	-0.0161	1.0000	0.0410	0.0102	0.0761	0.0593	0.0003	-0.0553	
maximum_nights	-0.0386	-0.0149	-0.0081	-0.0174	0.0410	1.0000	0.0223	0.0357	0.0520	0.0527	0.1483	
availability_30	-0.1092	-0.1176	-0.1127	-0.1100	0.1032	0.0223	1.0000	0.8864	0.8121	0.5228	-0.0991	
availability_60	-0.0689	-0.0704	-0.0609	-0.0602	0.0761	0.0357	0.8864	1.0000	0.9556	0.6491	-0.0727	
availability_90	-0.1040	-0.0929	-0.0912	-0.0906	0.0593	0.0520	0.8121	0.9556	1.0000	0.7150	-0.0423	
availability_365	-0.0749	-0.0515	-0.0386	-0.0517	0.0003	0.0527	0.5228	0.6491	0.7150	1.0000	0.0642	
number_of_reviews	-0.2031	-0.1713	-0.1196	-0.1590	-0.0553	0.1483	-0.0991	-0.0727	-0.0423	0.0642	1.0000	
number_of_reviews_itm	-0.2250	-0.2003	-0.1376	-0.1673	-0.0770	0.0968	-0.0836	-0.0405	0.0029	0.1176	0.8060	
number_of_reviews_l30d	-0.1697	-0.1600	-0.1008	-0.1253	-0.0587	0.0378	-0.0232	0.0261	0.0594	0.1224	0.5416	
review_scores_rating	0.0137	0.0519	0.0299	0.0082	-0.1125	-0.0468	-0.1841	-0.1747	-0.1770	-0.1551	0.0423	
review_scores_accuracy	0.0130	0.0471	0.0193	0.0015	-0.0283	-0.0334	-0.1848	-0.1729	-0.1747	-0.1357	0.0674	
review_scores_cleanliness	-0.0073	0.0253	0.0096	-0.0115	0.0097	-0.0261	-0.1389	-0.1221	-0.1296	-0.1291	0.0446	
review_scores_checkin	0.0024	0.0287	0.0161	-0.0058	0.0006	-0.0395	-0.1824	-0.1586	-0.1499	-0.1106	0.1068	
review_scores_communication	0.0182	0.0184	0.0303	0.0107	-0.0185	-0.0379	-0.2139	-0.1886	-0.1900	-0.1580	0.1083	
review_scores_location	0.0618	0.0756	0.0350	0.0411	-0.0151	-0.0340	-0.1357	-0.1356	-0.1467	-0.1322	-0.0337	
review_scores_value	-0.0227	-0.0054	-0.0239	-0.0328	0.0197	-0.0218	-0.1998	-0.2086	-0.2133	-0.1878	0.1046	
calculated_host_listings_count	0.2183	0.2211	0.1529	0.1976	-0.0765	0.0178	0.0937	0.0885	0.0973	0.1415	-0.2719	
reviews_per_month	-0.2214	-0.2367	-0.1442	-0.1741	-0.0617	-0.0054	-0.0592	-0.0187	0.0202	0.0892	0.6451	

b) How many factors do you get? Label at least the first three (If you get less than 3 factors, label all your factors)

To avoid multicollinearity, we conducted the PCA to summarize large number of similar/redundant variables into several factors. As a result, we got 6 components with their Eigenvalues > 1. These components capture approximately 75.93% (75.928%) of information contained in our 22 variables.



Renamed 6 Factors:

1. Ratings
2. Number of Reviews
3. Instant Availability
4. Room
5. Minimum Nights
6. Maximum Nights

Also, Multivariate data shows that these factors are not correlated.

Multivariate

Correlations

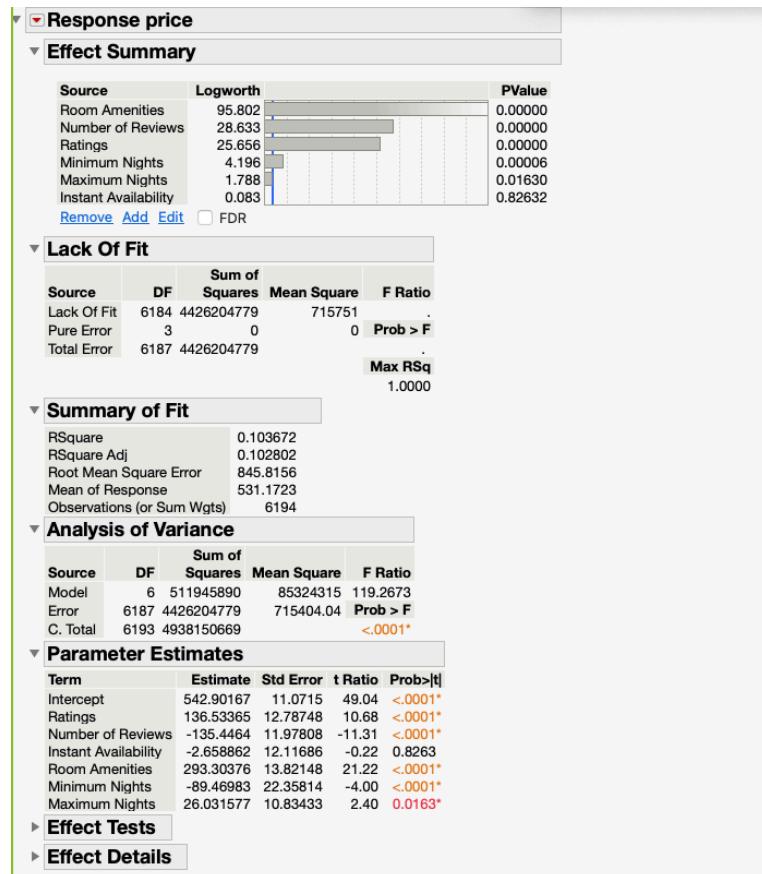
	Ratings	Number of Reviews	Instant Availability	Room Amenities	Minimum Nights	Maximum Nights
Ratings	1.0000	0.0027	-0.0298	-0.0098	0.1818	-0.0240
Number of Reviews	0.0027	1.0000	-0.0209	0.0333	0.1780	-0.0142
Instant Availability	-0.0298	-0.0209	1.0000	0.0224	-0.2551	-0.0148
Room Amenities	-0.0098	0.0333	0.0224	1.0000	-0.0396	-0.0071
Minimum Nights	0.1818	0.1780	-0.2551	-0.0396	1.0000	-0.1673
Maximum Nights	-0.0240	-0.0142	-0.0148	-0.0071	-0.1673	1.0000

c) Run a linear regression model on your y-variable using all the factors. What is the most important factor? If that factor was not labeled in Step b, give a name to that factor

From this linear regression, we can analyze that the most important factor is the Room, followed by Number of reviews, Ratings, Minimum nights and maximum nights.

Room amenities (# of bedrooms, beds, bathrooms, and accommodation), number of reviews and ratings do matter in predicting our y-variable, price per night.

When looking at the t-value, all values have probability of less than 0.05, except instant availability. In other words, these are all statistically significant with our y-variable, price per night. (Instant availability and maximum nights are not statistically significant in impacting the price per night).



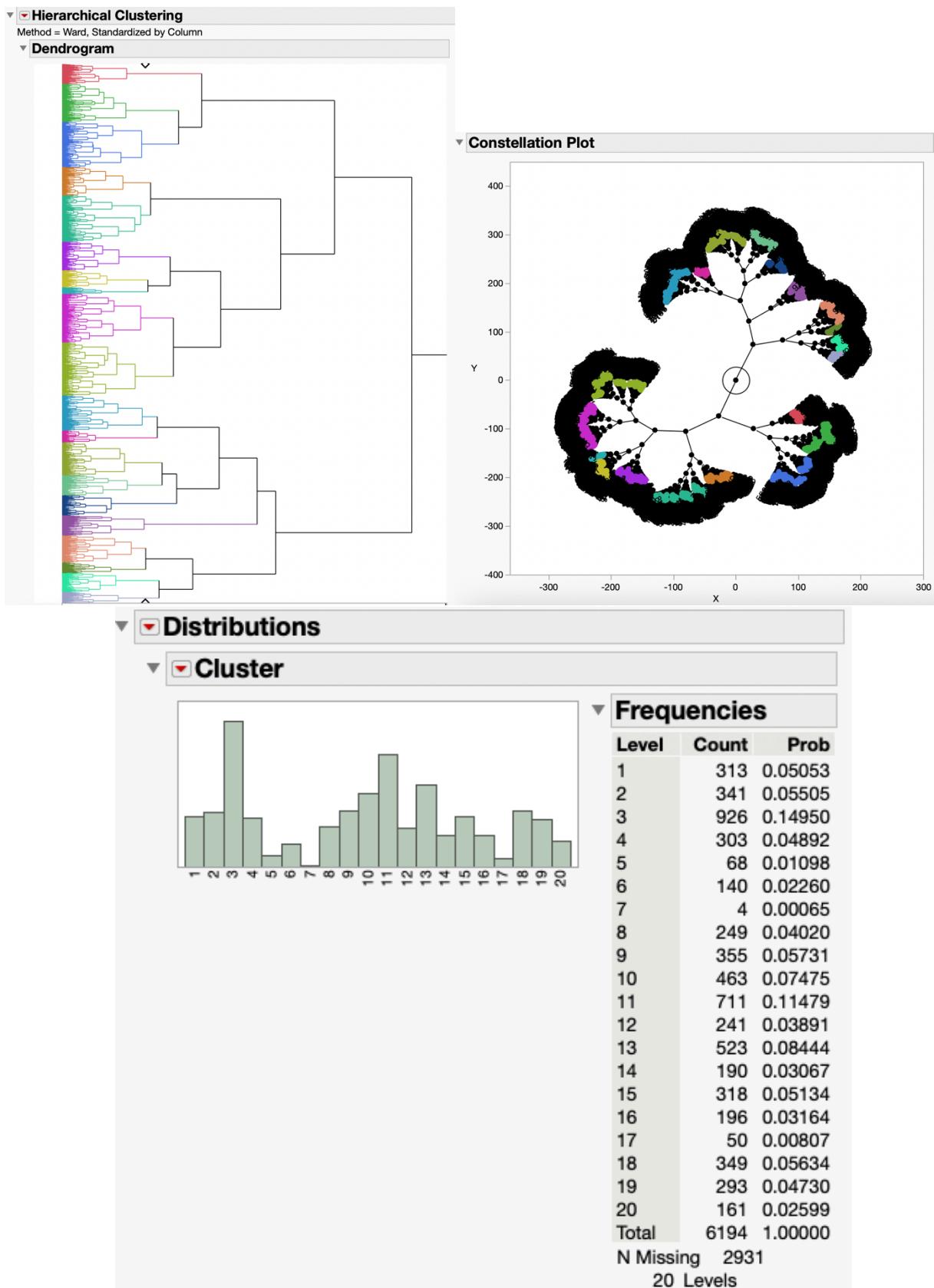
Clustering:

a) Run a Cluster Analysis using the factors from PCA/Factor Analysis AND other variables that you would like to add (variables that were not used when running PCA/Factor analysis). How many clusters are there? (Feel free to visualize the clusters using a Constellation Plot)

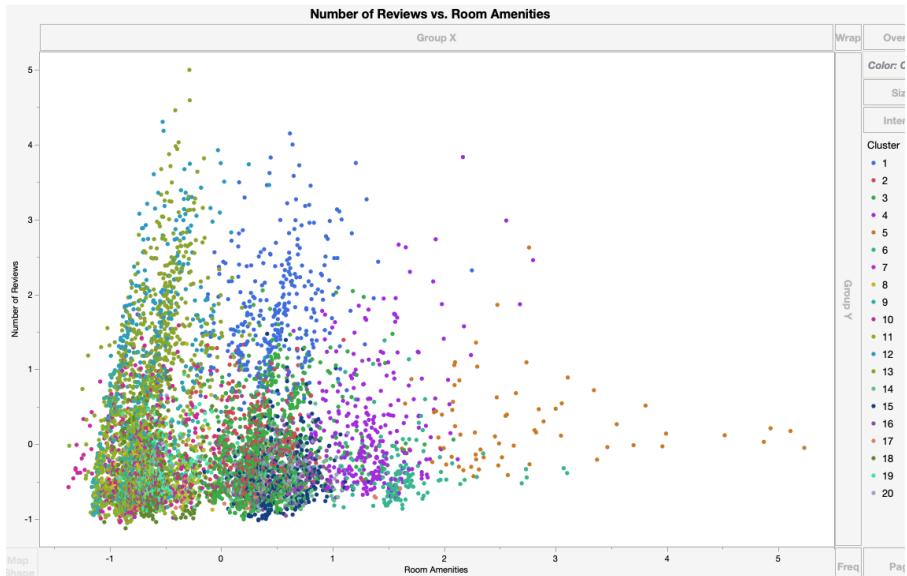
6 Factors + Room type + Neighborhood_cleansed + host_is_superhost

There are 20 clusters

Since we have many listings for each cluster, it created a condensed constellation plot.



b) Build a perceptual map using the two most important factors (the first two factors) and describe the patterns you detect.



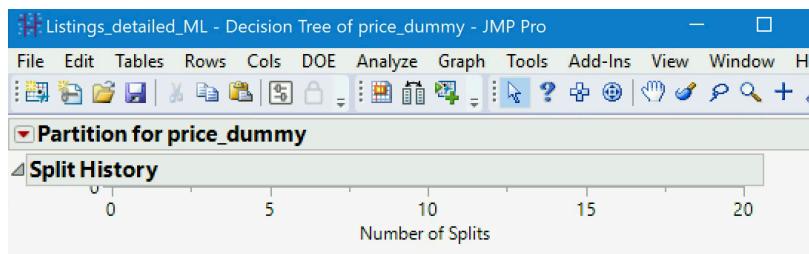
Each cluster (same colors) are in the same/similar area. Orange cluster tends to have higher number in room (# of bathrooms, bedrooms, beds, and accommodation).
Light blue cluster seems to have higher number of reviews.

Supervised Learning:

2 ML Models: Decision Tree, Bootstrapping Forest, and Neural Network

Y variable is price, 0 = below \$387 and 1 = above \$387

Decision Tree: 8 Splits

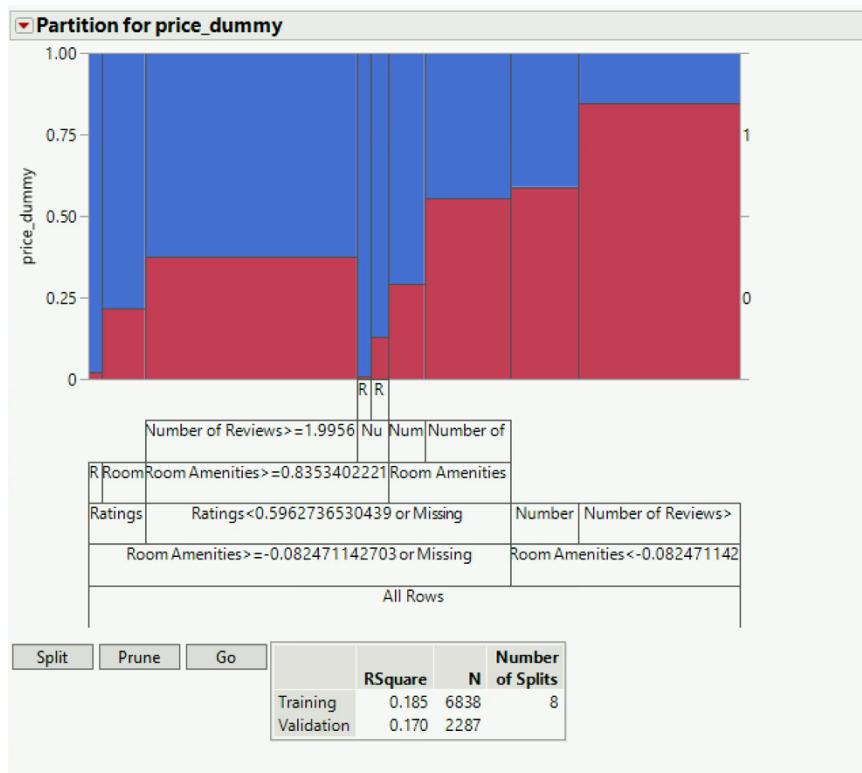


Fit Details

Measure	Training	Validation	Definition
Entropy RSquare	0.1853	0.1695	$1 - \text{Loglike}(\text{model})/\text{Loglike}(0)$
Generalized RSquare	0.3020	0.2792	$(1 - L(0)/L(\text{model}))^{(2/n)} / (1 - L(0)^{(2/n)})$
Mean -Log p	0.5647	0.5756	$\sum -\text{Log}(p[j])/n$
RASE	0.4391	0.4436	$\sqrt{\sum (y[j] - p[j])^2/n}$
Mean Abs Dev	0.3858	0.3903	$\sum y[j] - p[j] /n$
Misclassification Rate	0.2964	0.3039	$\sum (p[j] \neq p_{\text{Max}})/n$
N	6838	2287	n

Confusion Matrix

		Training		Validation	
		Actual	Predicted	Actual	Predicted
price_dummy	Count	0	1	0	1
	0	2350	1066	789	348
price_dummy	Rate	0	1	0	1
	0	0.688	0.312	0.694	0.306
price_dummy	Rate	1	0	1	0
	1	0.281	0.719	0.302	0.698



Bootstrapping Forest: 27.98% (0.2798)

Listings_detailed_ML - Bootstrap Forest of price_dummy - JM... - X

File Edit Tables Rows Cols DOE Analyze Graph Tools Add-Ins View
Window Help

Bootstrap Forest for price_dummy

Specifications

Target	price_dummy	Training Rows:	6838
Validation Column:	Validation	Validation Rows:	2287
		Test Rows:	0
Number of Trees in the Forest:	100	Number of Terms:	6
Number of Terms Sampled per Split:	5	Bootstrap Samples:	6838
		Minimum Splits per Tree:	10
		Minimum Size Split:	9

Overall Statistics

Measure	Training	Validation	Definition
Entropy RSquare	0.3892	0.2131	$1 - \text{Loglike}(\text{model})/\text{Loglike}(0)$
Generalized RSquare	0.5560	0.3410	$(1 - L(0)/L(\text{model}))^{(2/n)} / (1 - L(0)^{(2/n)})$
Mean -Log p	0.4234	0.5454	$\sum -\text{Log}(p[j])/n$
RASE	0.3694	0.4309	$\sqrt{\sum (y[j] - p[j])^2/n}$
Mean Abs Dev	0.3146	0.3729	$\sum y[j] - p[j] /n$
Misclassification Rate	0.1925	0.2798	$\sum (p[j] \neq p_{\text{Max}})/n$
N	6838	2287	n

Confusion Matrix

Training		Validation	
Actual	Predicted Count	Actual	Predicted Count
price_dummy	0 1	price_dummy	0 1
0	2418 998	0	726 411
1	318 3104	1	229 921

Actual	Predicted Rate		Actual	Predicted Rate	
price_dummy	0	1	price_dummy	0	1
0	0.708	0.292	0	0.639	0.361
1	0.093	0.907	1	0.199	0.801

Neural Network: 28.33% (0.2833)

Listings_detailed_ML - Neural of price_dummy - JMP Pro

File Edit Tables Rows Cols DOE Analyze Graph Tools Add-Ins View Window Help

Neural

Validation Column: Validation
Informative Missing

Model Launch

Model NTanH(3)

Training Validation

price_dummy		price_dummy	
Measures	Value	Measures	Value
Generalized RSquare	0.3425409	Generalized RSquare	0.3250907
Entropy RSquare	0.2141915	Entropy RSquare	0.2015995
RASE	0.4307987	RASE	0.4339492
Mean Abs Dev	0.3705252	Mean Abs Dev	0.3743613
Misclassification Rate	0.2860486	Misclassification Rate	0.2833406
-LogLikelihood	3724.5265	-LogLikelihood	1265.617
Sum Freq	6838	Sum Freq	2287

Confusion Matrix

		Actual	Predicted	Count
		price_dummy	0 1	
0	2118	1298	716	
1	658	2764	421	

Confusion Rates

		Actual	Predicted	Rate
		price_dummy	0 1	
0	0.620	0.380	0.630	
1	0.192	0.808	0.370	

Confusion Matrix

		Actual	Predicted	Count
		price_dummy	0 1	
0	227	923	716	
1	227	923	421	

Confusion Rates

		Actual	Predicted	Rate
		price_dummy	0 1	
0	0.197	0.803	0.370	
1	0.630	0.370	0.620	

Model Performance Comparison:

The lowest Misclassification is 27.98% and the highest AUC % is at 78.66%: The best model is **Bootstrap Forest**

Listings_detailed_ML - Model Comparison - JMP Pro

File Edit Tables Rows Cols DOE Analyze Graph Tools Add-Ins View Window Help

Model Comparison Validation=Training

Predictors

Measures of Fit for price_dummy

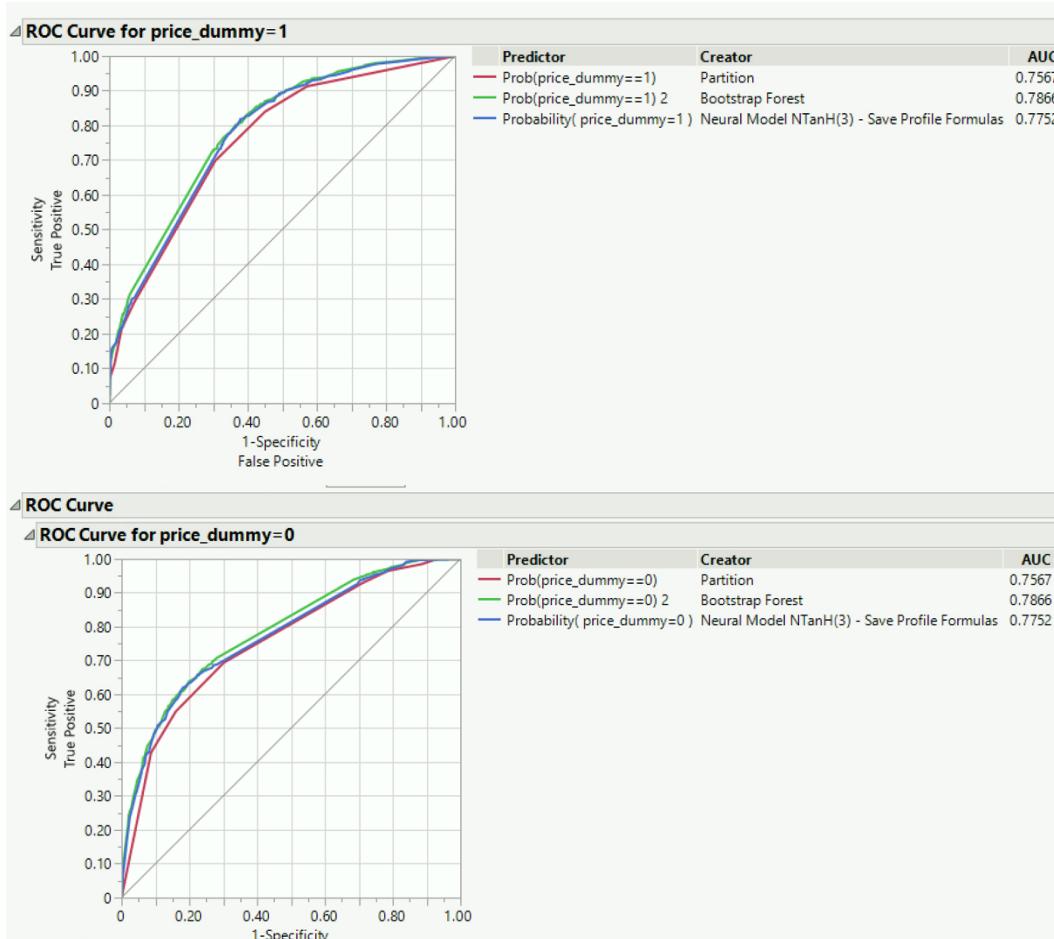
Creator	.2 .4 .6 .8	Entropy	Generalized	Mean	Misclassification	N		
Partition		RSquare	RSquare	-Log p	RASE	Abs Dev	Rate	
Bootstrap Forest		0.1853	0.3020	0.5647	0.4391	0.3858	0.2964	6838
Neural Model NTanH(3) - Save Profile Formulas		0.3892	0.5560	0.4234	0.3694	0.3146	0.1925	6838
		0.2142	0.3425	0.5447	0.4308	0.3705	0.2860	6838

Model Comparison Validation=Validation

Predictors

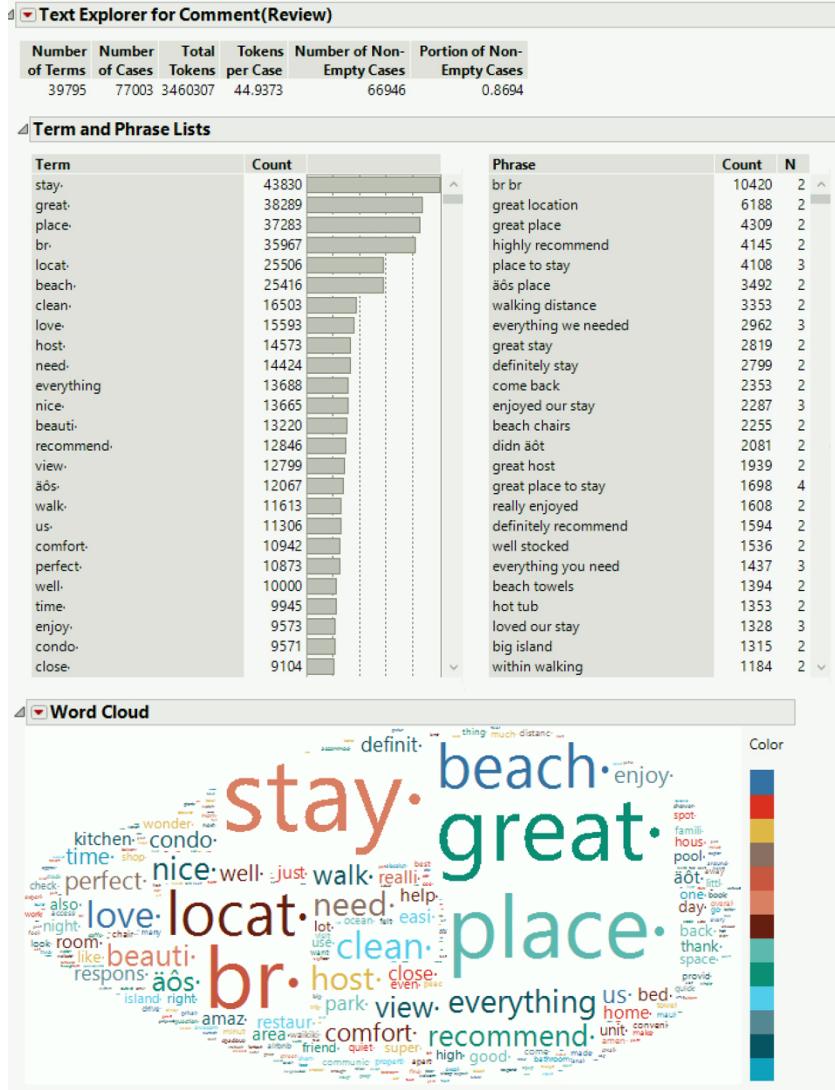
Measures of Fit for price_dummy

Creator	.2 .4 .6 .8	Entropy	Generalized	Mean	Misclassification	N		
Partition		RSquare	RSquare	-Log p	RASE	Abs Dev	Rate	
Bootstrap Forest		0.1695	0.2792	0.5756	0.4436	0.3903	0.3039	2287
Neural Model NTanH(3) - Save Profile Formulas		0.2131	0.3410	0.5454	0.4309	0.3729	0.2798	2287
		0.2016	0.3251	0.5534	0.4339	0.3744	0.2833	2287

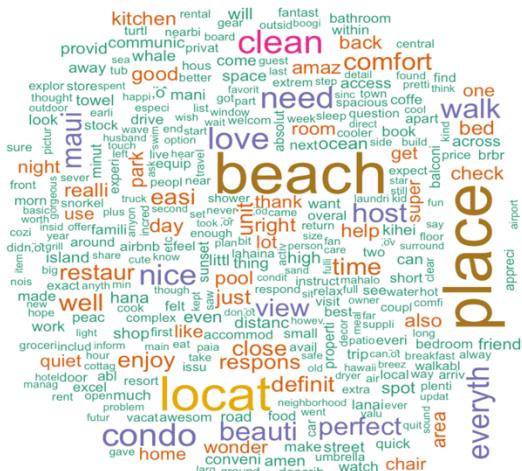


Text Mining: on (r and jmp)

Word Count: As shown below, most of the reviews consist of terms and phrases, such as “have stay, great, place, locat, beach, clean, love, etc. (JMP had several unrecognizable letters (foreign letters), but we were able to generate better word cloud with R using cleaned up data).

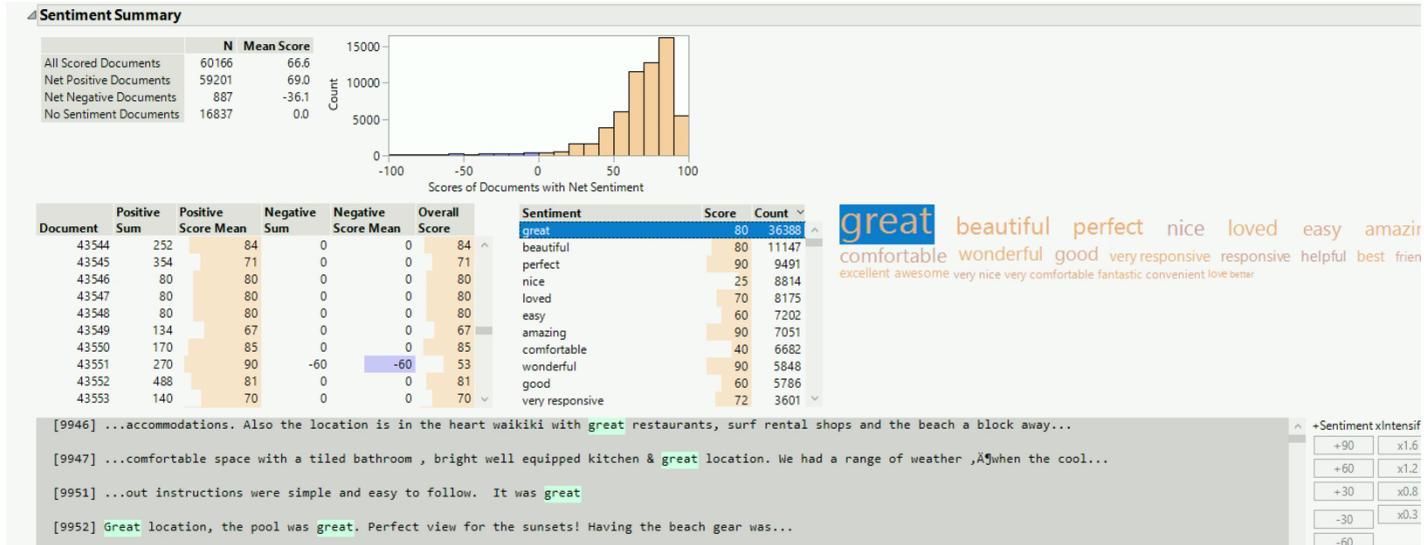


R- with cleaned data



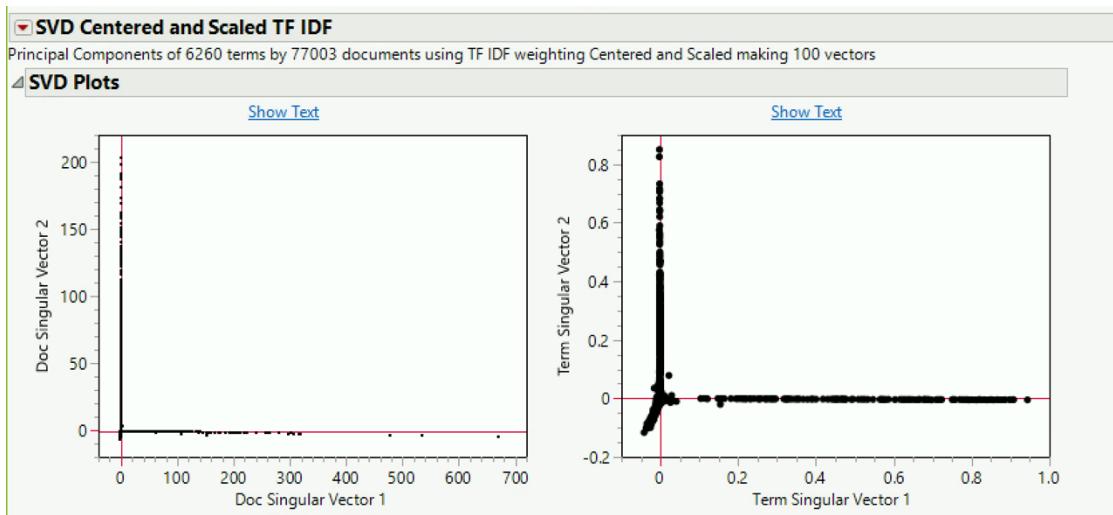
Sentiment Analysis

As we explained previously, these terms and phrases are very positive, it is proven by the sentiment analysis that almost 70% are positive, including the word, “great, beautiful, perfect, loved, easy, amazing, comfortable, wonderful, etc.”



Topic Modeling

With the topic modeling, Topic 1 to 5 has many encoded symbols where the system cannot translate into English. This means that the reviews are left with other foreign languages. (Topic 2 contains German and Topic 3 contains French, Topic 5 and 9 contains Spanish or Portuguese). Topic 6 contains amenities, and Topic 8 contains things like location and amenities and nearby activities).



Top Loadings by Topic

Topic 1		Topic 2		Topic 3		Topic 4		Topic 5	
Term	Loading								
éæ.	0.94361	und-	0.86780	et.	0.85140	µ-	0.99052	y-	0.81039
éø.	0.90982	die-	0.84462	le-	0.76915	a-	0.98894	el-	0.80161
éß.	0.90733	der-	0.74897	est-	0.72907	å-	0.98840	muy-	0.70790
éü.	0.90537	ist-	0.73139	pour-	0.72051	æ-	0.98492	todo-	0.67102
éñ.	0.90040	das-	0.72653	v+	0.71535	ω-	0.98100	lo-	0.64491
éó.	0.89744	ein-	0.72034	de-	0.71488	≤-	0.98095	una-	0.64312
éí.	0.88591	war-	0.69604	nous-	0.68387	¥-	0.97814	para-	0.60221
éää.	0.87830	sehr-	0.69481	les-	0.67725	ç-	0.97698	las-	0.59856
éń.	0.87449	wir-	0.68151	trv. ®-	0.67159	Π-	0.97510	en-	0.58730
éä.	0.87316	mit-	0.65893	pas-	0.64631	è-	0.97309	que-	0.57554
éç.	0.86564	zu-	0.65506	il-	0.60831	ã-	0.97292	la-	0.57048
éȝ.	0.85628	nicht-	0.65389	des-	0.60202	∞-	0.97223	tien-	0.56610
éää.	0.85070	fvør-	0.63303	tout-	0.59736	ʃ-	0.94761	playa-	0.54976
éç.	0.84624	auch-	0.60403	avon-	0.59070	á-	0.94656	lugar-	0.52945
çä.	0.83951	man-	0.59288	la-	0.58836	ø-	0.94528	pero-	0.51864
éø.	0.83923	auf-	0.57522	avec-	0.56940	º-	0.94065	los-	0.50595
éö.	0.83152					π-	0.93717		
éé.	0.82804					ð-	0.93581		
						é-	0.92873		

Topic 6		Topic 7		Topic 8		Topic 9		Topic 10	
Term	Loading	Term	Loading	Term	Loading	Term	Loading	Term	Loading
not	0.53772	ໄຕ້.	0.61930	beach-	0.64610	o-	0.80424	us-	0.31623
no	0.40702	ໄຕ້ເນື້ອທີ່.	0.50575	walk-	0.44119	✓@-	0.78416	fruit-	0.30630
a&t-	0.35682	ໄຕ້ເອົາເກີ-.	0.50261	chair-	0.43263	uma-	0.75327	home-	0.28915
br-	0.33159	ໄຕ້ເບົ-	0.48032	restaur-	0.37231	e-	0.74904	so	0.27440
only	0.33086	ໄຕ້ວິນ-	0.47011	kitchen-	0.36897	muito-	0.74262	island-	0.27288
dirti-	0.32102	ໄຕ້ວັດ-.	0.44988	pool-	0.34977	n'£o-	0.69338	stay-	0.27086
door-	0.31944	ໄຕ້ດ-	0.44599	towel-	0.34203	bem-	0.68972	tree-	0.26307
bathroom	0.29513	ໄຕ້ວິນ-	0.44222	condo-	0.32185	com-	0.68753	frog-	0.25603
one-	0.29179	ໄຕ້ວິນເລື່ອມຫຼັກສິດ-	0.42943	shop-	0.32133	foi-	0.67559	peac-	0.24899
issu-	0.28791	ໄຕ້ວິນເສົ້າ-	0.42748	well-	0.31586	mas-	0.65521	feel-	0.24709
if-	0.28539	ໄຕ້ວິນເລື່ອມຫຼັກສິດ-	0.42060	board-	0.31524	localizav@v£o-	0.60706	beautiful	0.24512
room-	0.27292	ໄຕ້ວິນເລື່ອມຫຼັກສິດ-	0.42059	snorkel-	0.31130	✓@tima-	0.60190	bird-	0.23804
stain-	0.27238	ໄຕ້ວິນ-	0.41724	locat-	0.30866	poi-	0.60001	hous-	0.23602
so	0.27227	ໄຕ້ວິນ-	0.41698	umbrella-	0.29714	praia-	0.53940	fresh-	0.23307
shower-	0.26954	ໄຕ້ວິນໂຄ-	0.41317	also-	0.28930			recommend-	0.22724
use-	0.26923	ໄຕ້ວິນ-	0.40824					experi-	0.22597
		ໄນ່.	0.39699					natur-	0.22194

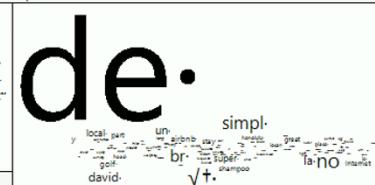
Topic 7



Topic 8



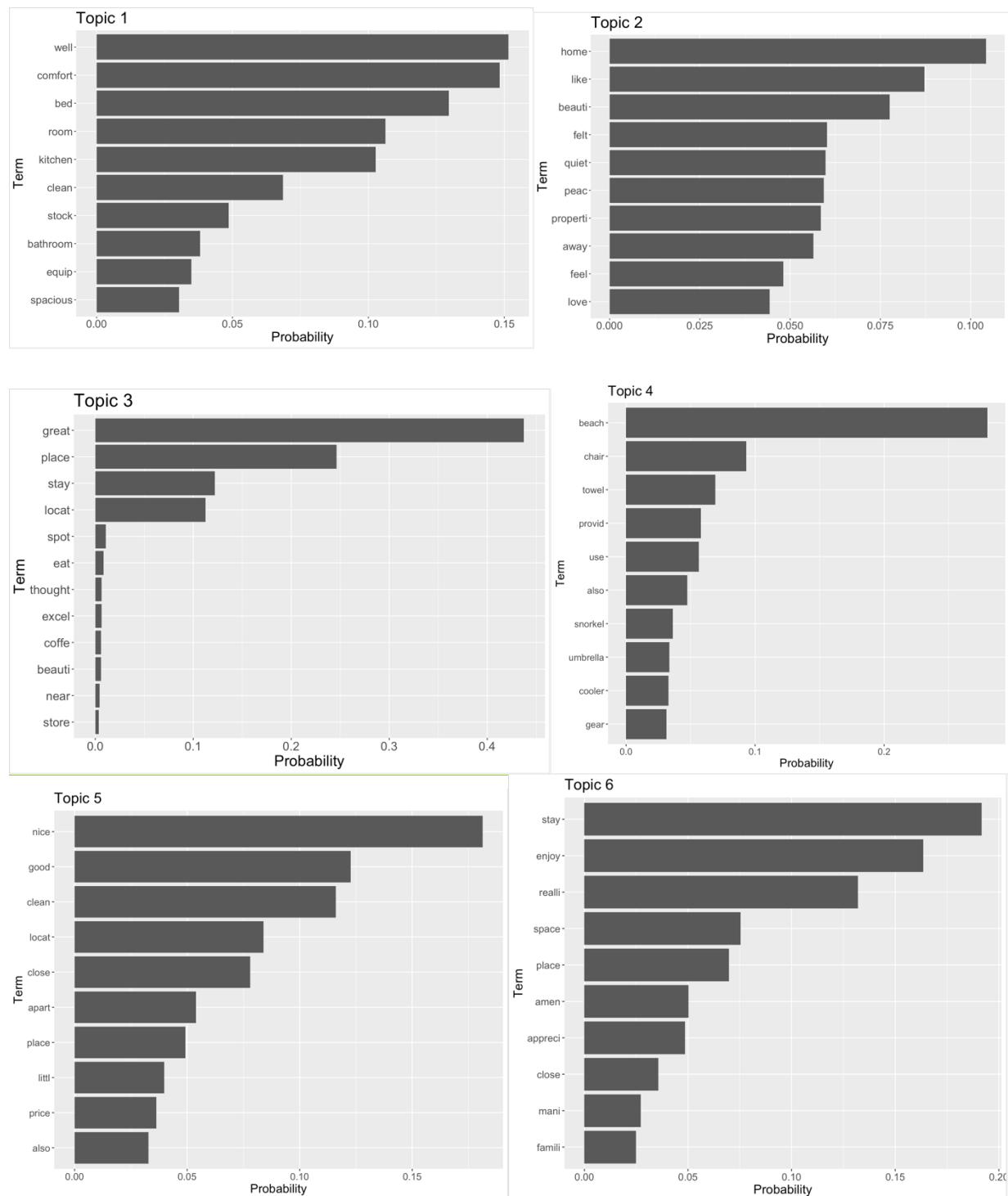
Topic 9

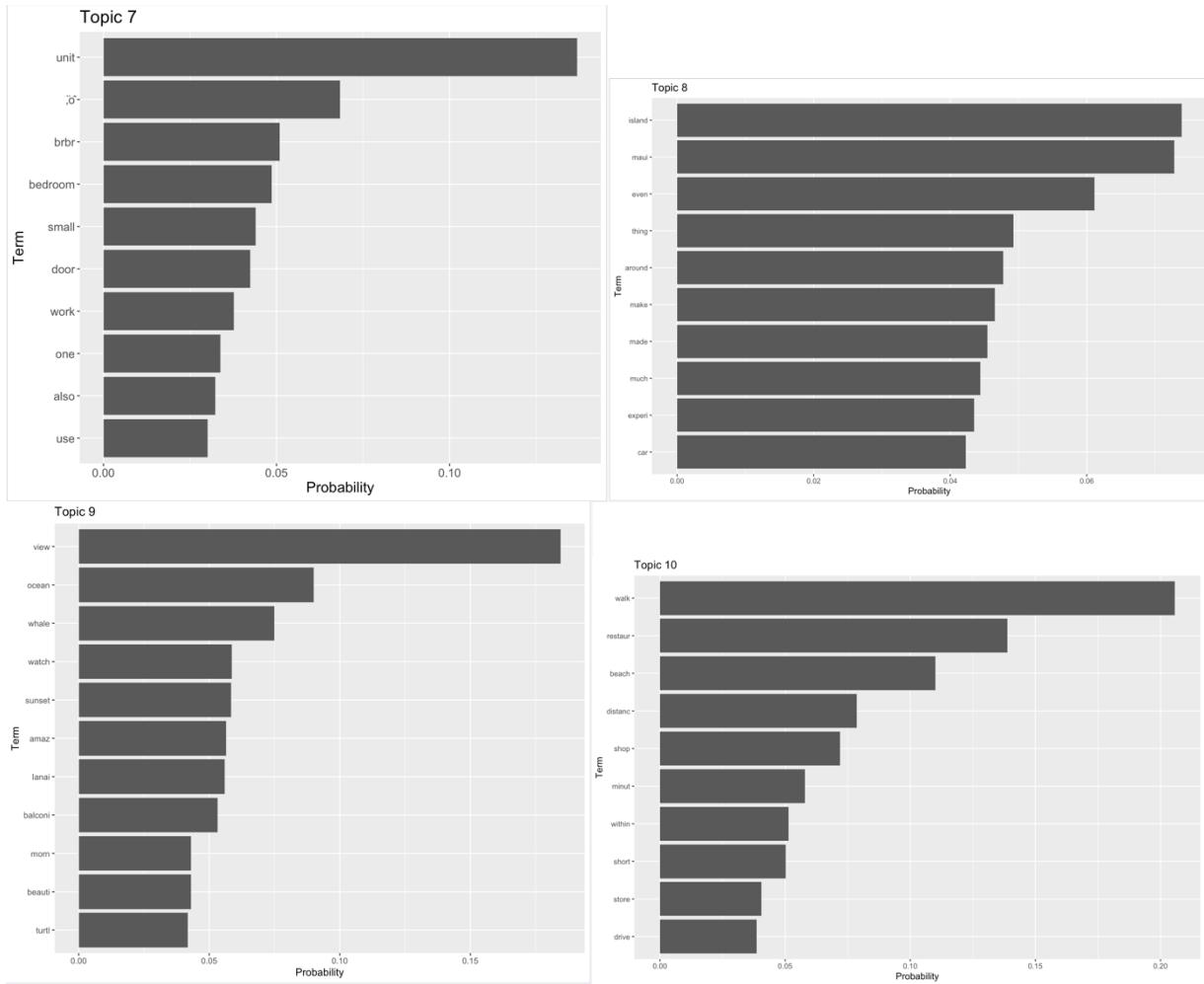


Topic 10



With cleaned data, we were able to generate better topic modeling with R



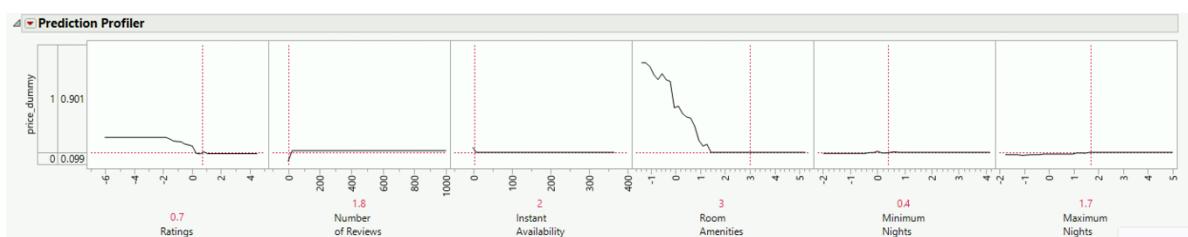


Price Prediction Profiler Comparison with Mid-term model

Bootstrap Forest Model – all 3 ML Scenarios

Price 0 = Below \$387, Price 1 = Above \$387

Scenario 1: More than 90% of our data is above \$387 when we have a good rating, quite a number of reviews, have better room amenities (# of bathrooms, bedrooms, and beds), has many availabilities, and have higher minimum nights and higher maximum nights.



Scenario 2: More than 60% of data has listings below \$387 when it has the lower ratings, number of reviews, less availability, room amenities (# of bathrooms, bedrooms, and beds), lower minimum nights, and lower maximum nights.

