

Advances in Data sciences

Final Project Report

<http://hoteladvisorsystem.azurewebsites.net/>

Team 3: Amitha Murali, Divyansh Srivastava, Jyoti Sharma

Table of Contents

1.1 Problem statement	3
1.2 Data Visualization.....	9
1.2.1 Power BI Analysis:.....	9
1.2.1 Tableau Analysis:.....	13
1.3 Data Pre-Processing.....	16
1.3.1 Missing Data Treatment.....	17
1.3.2 Feature Selection.....	18
1.3.3 Outliers detection and removal.....	19
1.4 Classification Models.....	22
1.4.1Overall Design	22
1.4.2 Azure models	23
1.4.3 Multiclass Decision Forest.....	24
1.4.4 One – vs – All Multiclass	26
1.4.4.1 Two class Boosted Decision Tree	27
1.4.4.2 Two class Decision Forest.....	29
1.4.5 Classification model comparison	31
1.4.6 Web Service	33
1.5 Recommendation System	36
1.5.1 Overall Design	36
1.5.2 Content Based Recommendation System	37
1.5.3 Hybrid Based Recommendation System.....	39
1.6 Sentimental Analysis.....	43
1.6.1Overall Design	44
1.6.2 Integration with Web Application.....	47
1.7 Web Application.....	47
1.8 References	54

Hotel Recommendation System

1.1 Problem statement

Expedia has provided us logs of customer behavior. These include what customers searched for, how they interacted with search results (click/book), whether or not the search result was a travel package. The data is a random selection from Expedia and is not representative of the overall statistics.

We are interested in predicting which hotel group a user is going to book. Expedia has in-house algorithms to form hotel clusters, where similar hotels for a search (based on historical price, customer star ratings, geographical locations relative to city center, etc) are grouped together. These hotel clusters serve as good identifiers to which types of hotels people are going to book, while avoiding outliers such as new hotels that don't have historical data.

Our goal is to predict the booking outcome (hotel cluster) for a user event, based on their search and other attributes associated with that user event.

Recommending hotels to users based on hotel properties and user behavior. Sentiment Analysis of the reviews given by the user.

The train and test datasets are split based on time: training data from 2013 and 2014, while test data are from 2015. Training data includes all the users in the logs, including both click events and booking events. We only use training dataset.

destinations.csv data consists of features extracted from hotel reviews text.

train.csv

Column name	Description	Data type
date_time	Timestamp	string
site_name	ID of the Expedia point of sale (i.e. Expedia.com, Expedia.co.uk, Expedia.co.jp, ...)	int
posa_continent	ID of continent associated with site_name	int
user_location_country	The ID of the country the customer is located	int

user_location_region	The ID of the region the customer is located	int
user_location_city	The ID of the city the customer is located	int
orig_destination_distance	Physical distance between a hotel and a customer at the time of search. A null means the distance could not be calculated	double
user_id	ID of user	int
is_mobile	1 when a user connected from a mobile device, 0 otherwise	tinyint
is_package	1 if the click/booking was generated as a part of a package (i.e. combined with a flight), 0 otherwise	int
Channel	ID of a marketing channel	int
srch_ci	Checkin date	string
srch_co	Checkout date	string
srch_adults_cnt	The number of adults specified in the hotel room	int
srch_children_cnt	The number of (extra occupancy) children specified in the hotel room	int
srch_rm_cnt	The number of hotel rooms specified in the search	int
srch_destination_id	ID of the destination where the hotel search was performed	int
srch_destination_type_id	Type of destination	int
hotel_continent	Hotel continent	int
hotel_country	Hotel country	int
hotel_market	Hotel market	int
is_booking	1 if a booking, 0 if a click	tinyint
Cnt	Numer of similar events in the context of the same user session	bigint
hotel_cluster	ID of a hotel cluster	int

destinations.csv

Column name	Description	Data type
srch_destination_id	ID of the destination where the hotel search was performed	int

Column name	Description	Data type
d1-d149	latent description of search regions	double

“Hotel” refers to hotels, apartments, B&Bs, hostels and other properties appearing on Expedia’s websites. Room types are not distinguished and the data can be assumed to apply to the least expensive room type.

Most of the data are for searches that resulted in a purchase, but a small proportion are for searches not leading to a purchase.

Column Name	Data Type	Description
srch_id	Integer	The ID of the search
date_time	Date/time	Date and time of the search
site_id	Integer	ID of the Expedia point of sale (i.e. Expedia.com, Expedia.co.uk, Expedia.co.jp,)
visitor_location_country_id	Integer	The ID of the country the customer is located
visitor_hist_starrating	Float	The mean star rating of hotels the customer has previously purchased; null signifies there is no purchase history on the customer
visitor_hist_adr_usd	Float	The mean price per night (in US\$) of the hotels the customer has previously purchased; null signifies there is no purchase history on the customer
prop_country_id	Integer	The ID of the country the hotel is located in
prop_id	Integer	The ID of the hotel
prop_starrating	Integer	The star rating of the hotel, from 1 to 5, in increments of 1. A 0 indicates the property has no stars, the star rating is not known or cannot be publicized.
prop_review_score	Float	The mean customer review score for the hotel on a scale out of 5, rounded to 0.5 increments. A 0 means there have been no reviews, null that the information is not available.
prop_brand_bool	Integer	+1 if the hotel is part of a major hotel chain; 0 if it is an independent hotel

prop_location_score1	Float	A (first) score outlining the desirability of a hotel's location
prop_location_score2	Float	A (second) score outlining the desirability of the hotel's location
prop_log_historical_price position	Float Integer	The logarithm of the mean price of the hotel over the last trading period. A 0 will occur if the hotel was not sold in that period. Hotel position on Expedia's search results page. This is only provided for the training data, but not the test data.
price_usd	Float	Displayed price of the hotel for the given search. Note that different countries have different conventions regarding displaying taxes and fees and the value may be per night or for the whole stay
promotion_flag	Integer	+1 if the hotel had a sale price promotion specifically displayed
gross_booking_usd	Float	Total value of the transaction. This can differ from the price_usd due to taxes, fees, conventions on multiple day bookings and purchase of a room type other than the one shown in the search
srch_destination_id	Integer	ID of the destination where the hotel search was performed
srch_length_of_stay	Integer	Number of nights stay that was searched
srch_booking_window	Integer	Number of days in the future the hotel stay started from the search date
srch_adults_count	Integer	The number of adults specified in the hotel room
srch_children_count	Integer	The number of (extra occupancy) children specified in the hotel room
srch_room_count	Integer	Number of hotel rooms specified in the search
srch_saturday_night_bool	Boolean	+1 if the stay includes a Saturday night, starts from Thursday with a length of stay is less than or equal to 4 nights (i.e. weekend); otherwise 0
srch_query_affinity_score	Float	The log of the probability a hotel will be clicked on in Internet searches (hence the values are negative) A null signifies there are

		no data (i.e. hotel did not register in any searches)
orig_destination_distance	Float	Physical distance between the hotel and the customer at the time of search. A null means the distance could not be calculated.
random_bool	Boolean	+1 when the displayed sort was random, 0 when the normal sort order was displayed
comp1_rate	Integer	+1 if Expedia has a lower price than competitor 1 for the hotel; 0 if the same; -1 if Expedia's price is higher than competitor 1; null signifies there is no competitive data
comp1_inv	Integer	+1 if competitor 1 does not have availability in the hotel; 0 if both Expedia and competitor 1 have availability; null signifies there is no competitive data
comp1_rate_percent_diff	Float	The absolute percentage difference (if one exists) between Expedia and competitor 1's price (Expedia's price the denominator); null signifies there is no competitive data
comp2_rate		
comp2_inv		(same, for competitor 2 through 8)
comp2_rate_percent_diff		
comp3_rate		
comp3_inv		
comp3_rate_percent_diff		
comp4_rate		
comp4_inv		
comp4_rate_percent_diff		
comp5_rate		
comp5_inv		
comp5_rate_percent_diff		
comp6_rate		
comp6_inv		
comp6_rate_percent_diff		
comp7_rate		
comp7_inv		
Comp7_rate_percent_diff		

Comp8_rate	
Comp8_inv	
Comp8_rate_percent_diff	

site_id

PLAN YOUR TRIP ON EXPEDIA

Flight Hotel Car Activities Cruise

Flight + Hotel Flight + Car Flight + Hotel + Car Hotel + Car

CHOOSE FROM MORE THAN 140,000 HOTELS WORLDWIDE

Hotel

Find hotels near:
A city, airport or attraction

What City?
New York (and vicinity), New York, United States of America

Check-in: 10/18/2013 Check-out: 10/20/2013 Rooms: 1

srch_destination_id

srch_room_count

srch_booking_window

Room 1 2 0

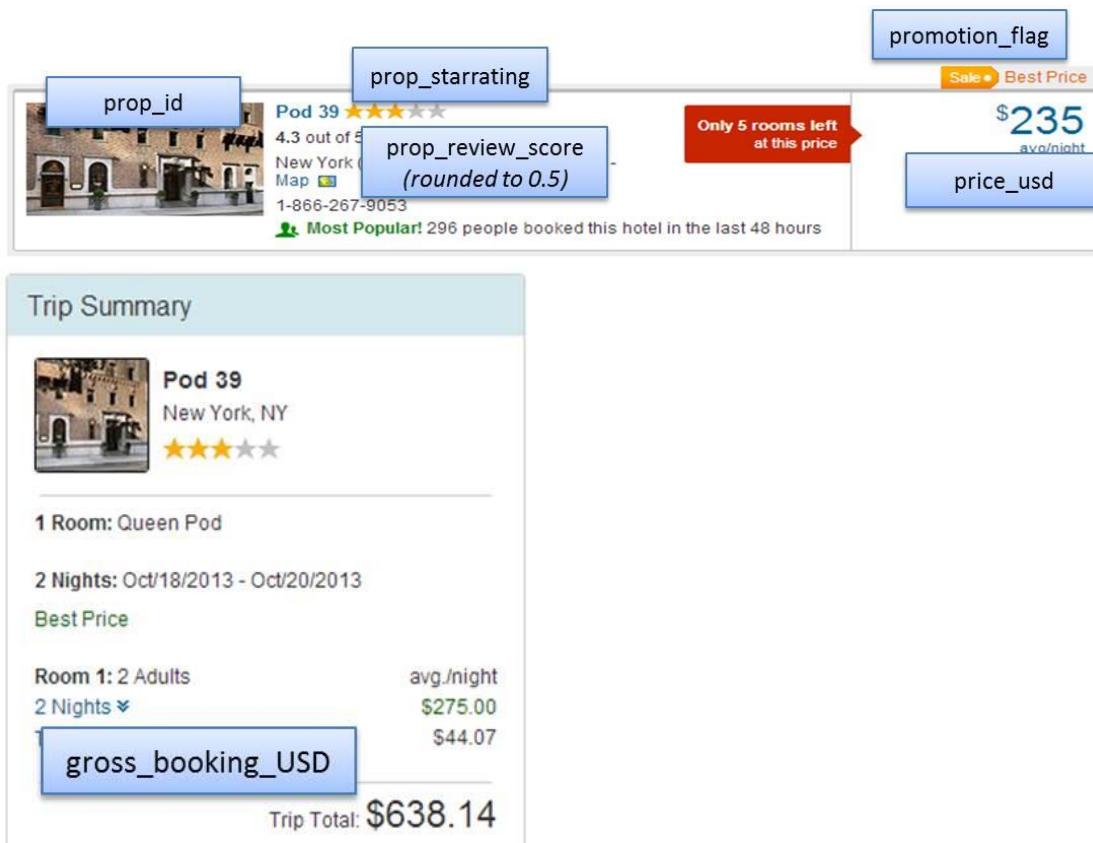
srch_length_of_stay

srch_adults_count

srch_children_count

BEST PRICE GUARANTEE

SEARCH FOR HOTELS



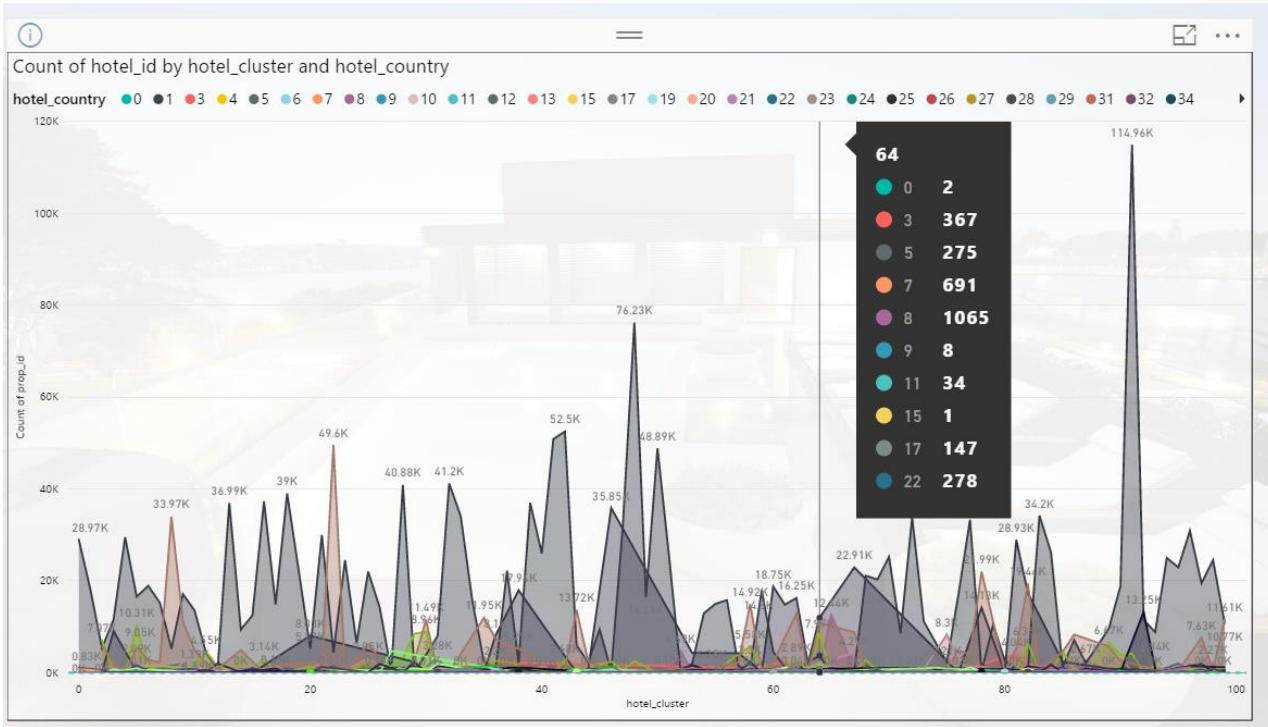
1.2 Data Visualization

Patterns, trends and correlations that might go undetected in text-based data can be exposed and recognized easier with data visualization software.

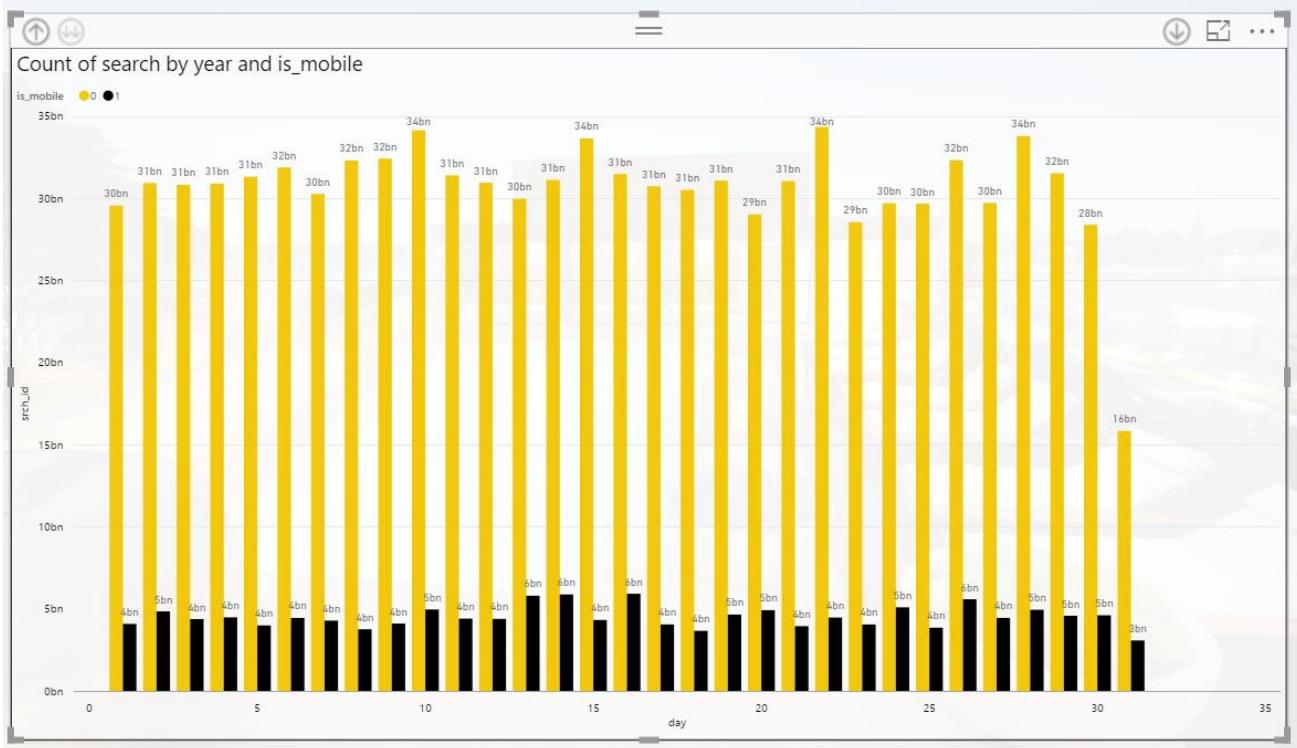
1.2.1 Power BI Analysis:

We have used the cleansed data to visualize and draw patterns about the hotel cluster data.

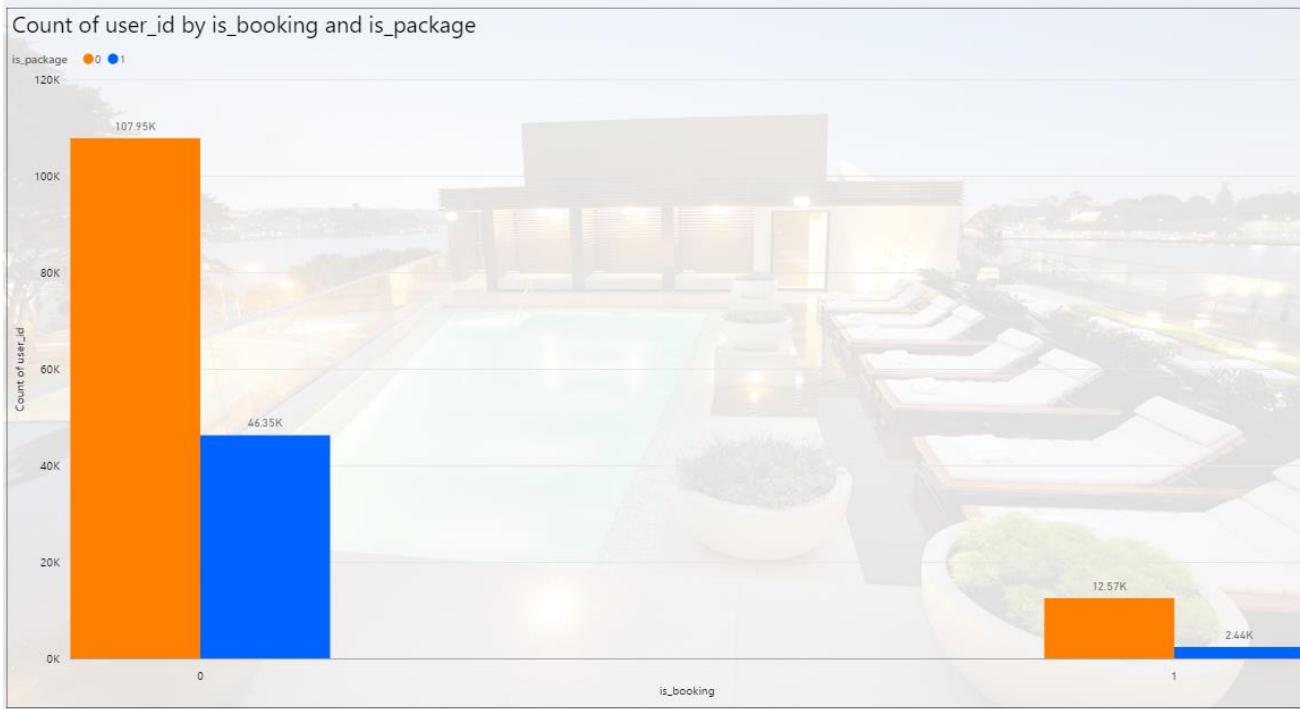
Graph1 – Count of hotel id by hotel cluster and hotel country



Graph2 – Count of search requests by year, month, day and is mobile factor.



Graph3 – Count of user id by booking and package.



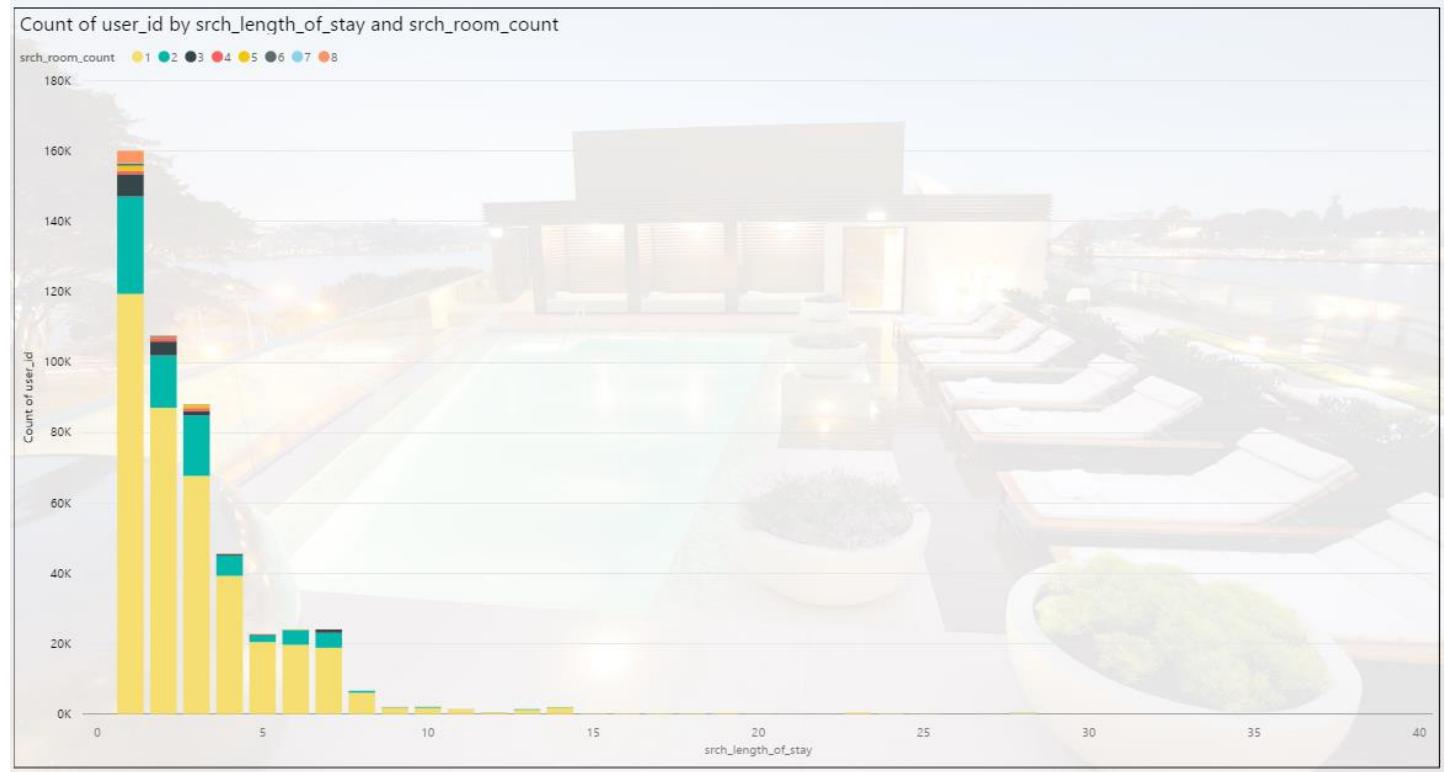
Graph4 – Average of gross bookings by hotel cluster.



Graph5 – Count of hotel id by hotel star rating and hotel brand bool



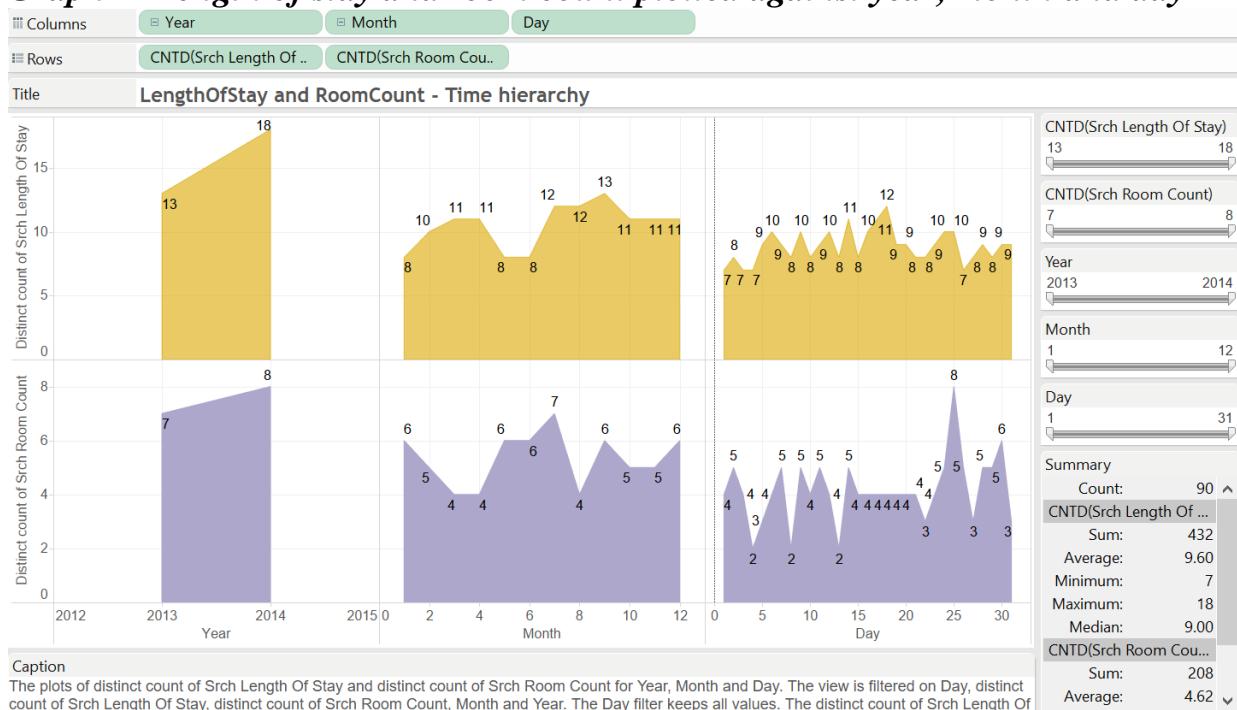
Graph6 – Count of user id by search length of stay and search room count



1.2.1 Tableau Analysis:

We have used the cleansed data to visualize and draw patterns about the hotel cluster data.

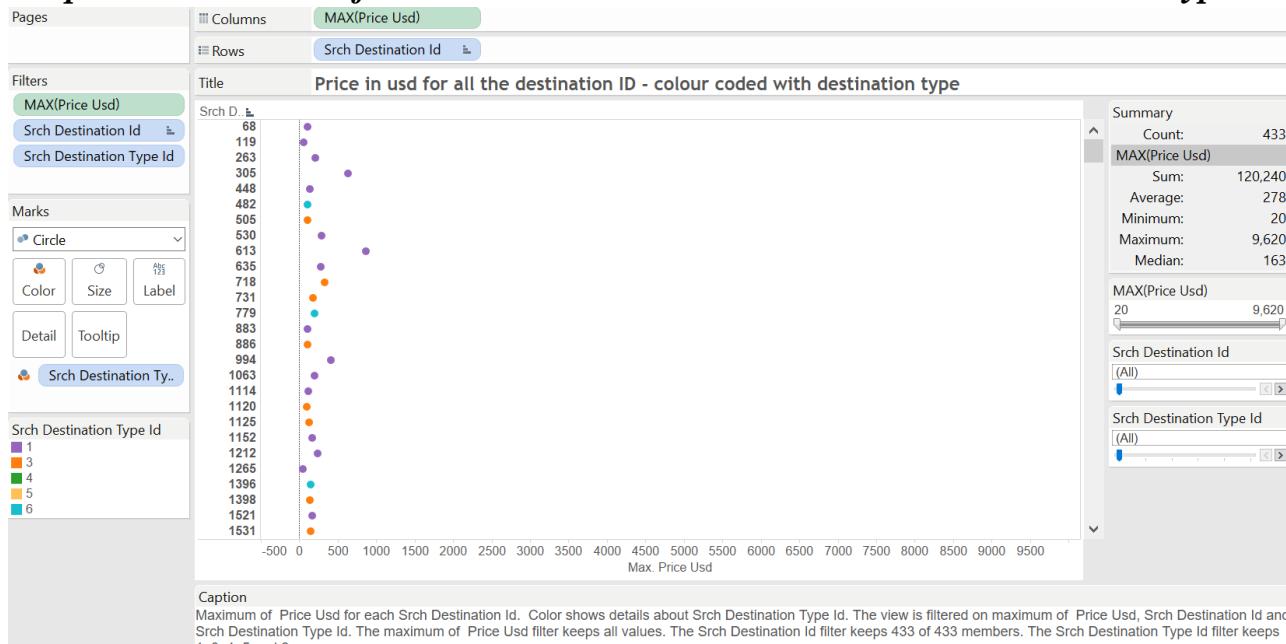
Graph1 –Length of stay and room count plotted against year, month and day



Caption

The plots of distinct count of Srch Length Of Stay and distinct count of Srch Room Count for Year, Month and Day. The view is filtered on Day, distinct count of Srch Length Of Stay, distinct count of Srch Room Count, Month and Year. The Day filter keeps all values. The distinct count of Srch Length Of

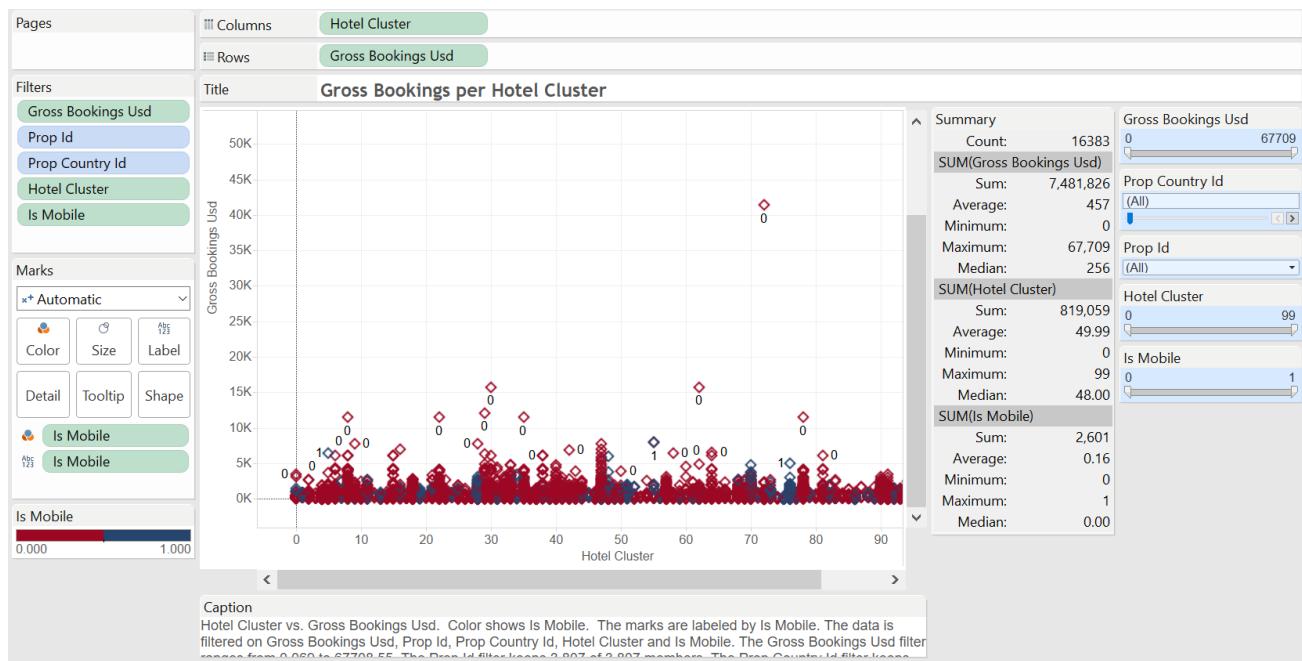
Graph2 – Price in usd for all destination Id – colour coded with destination type



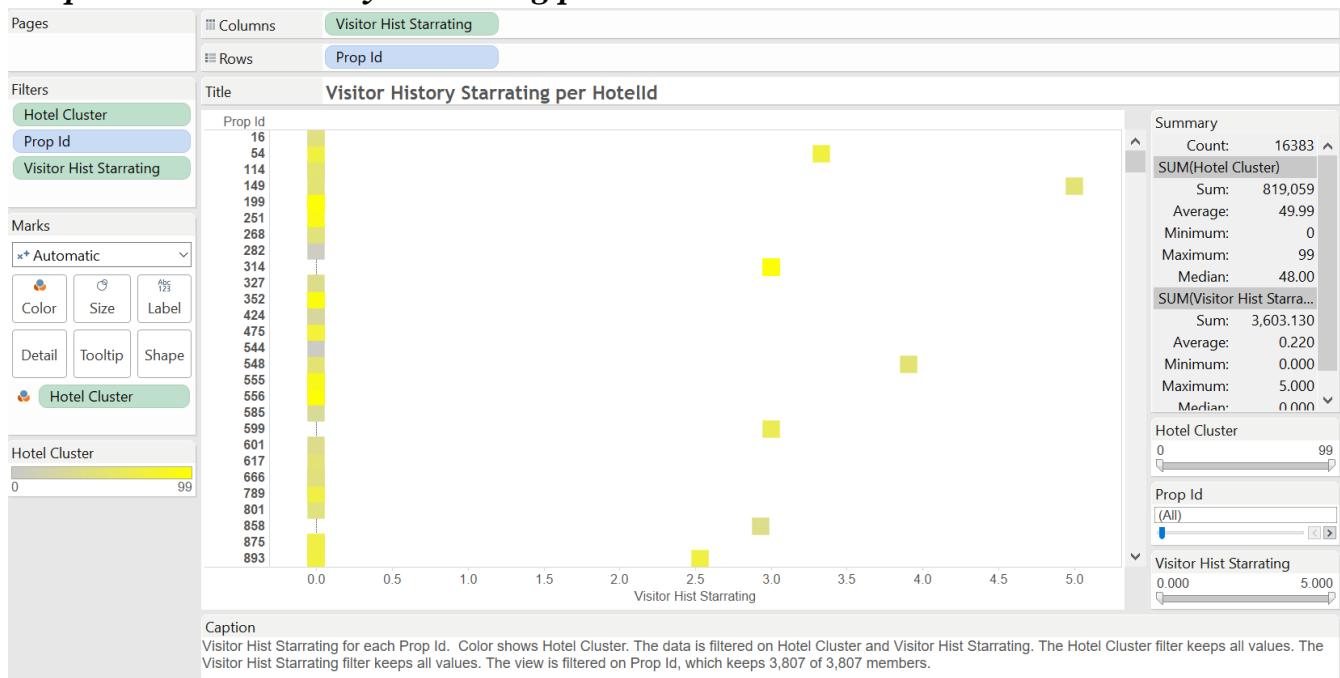
Caption

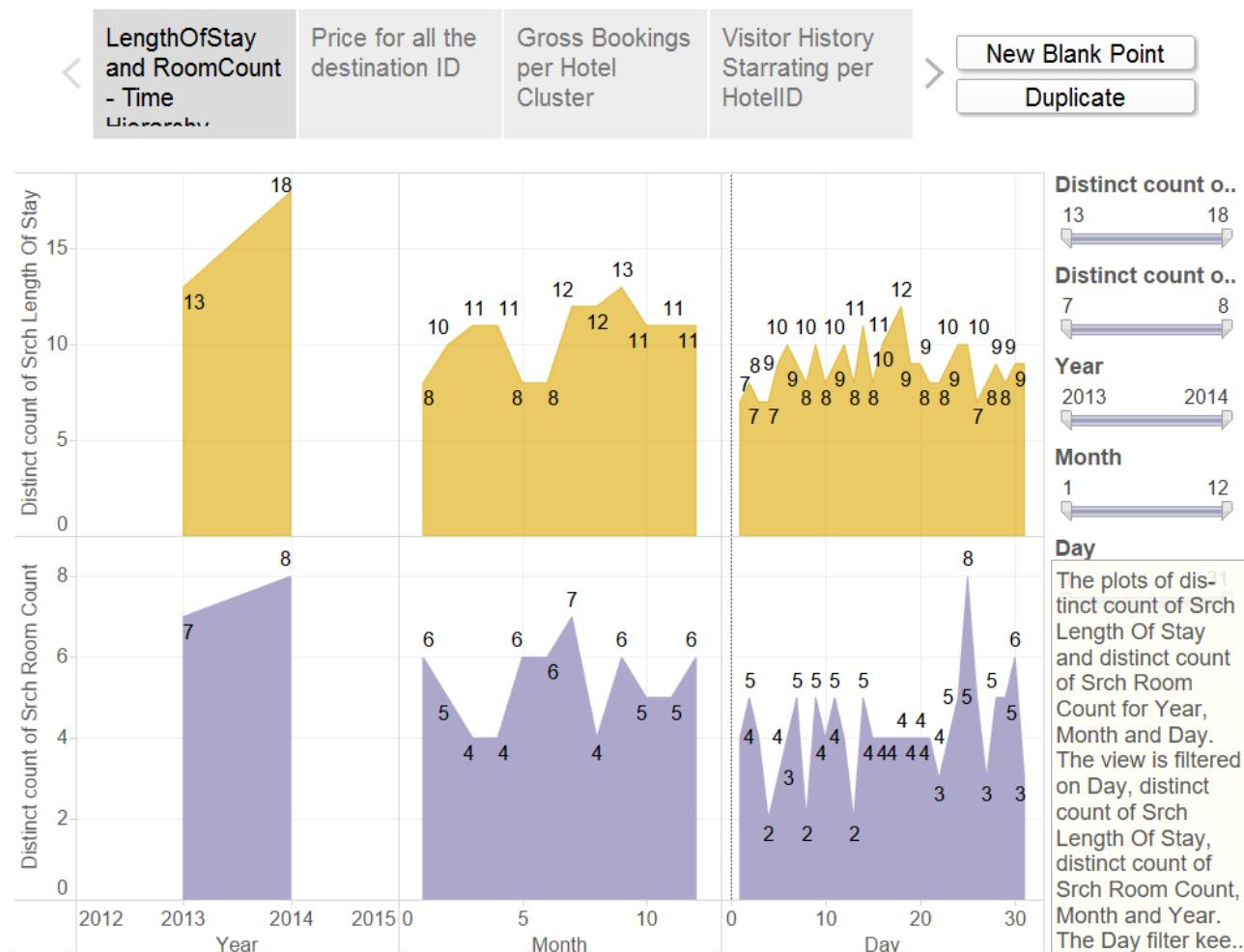
Maximum of Price Usd for each Srch Destination Id. Color shows details about Srch Destination Type Id. The view is filtered on maximum of Price Usd, Srch Destination Id and Srch Destination Type Id. The maximum of Price Usd filter keeps all values. The Srch Destination Id filter keeps 433 of 433 members. The Srch Destination Type Id filter keeps 6 of 6 members.

Graph3–Gross bookings per hotel cluster



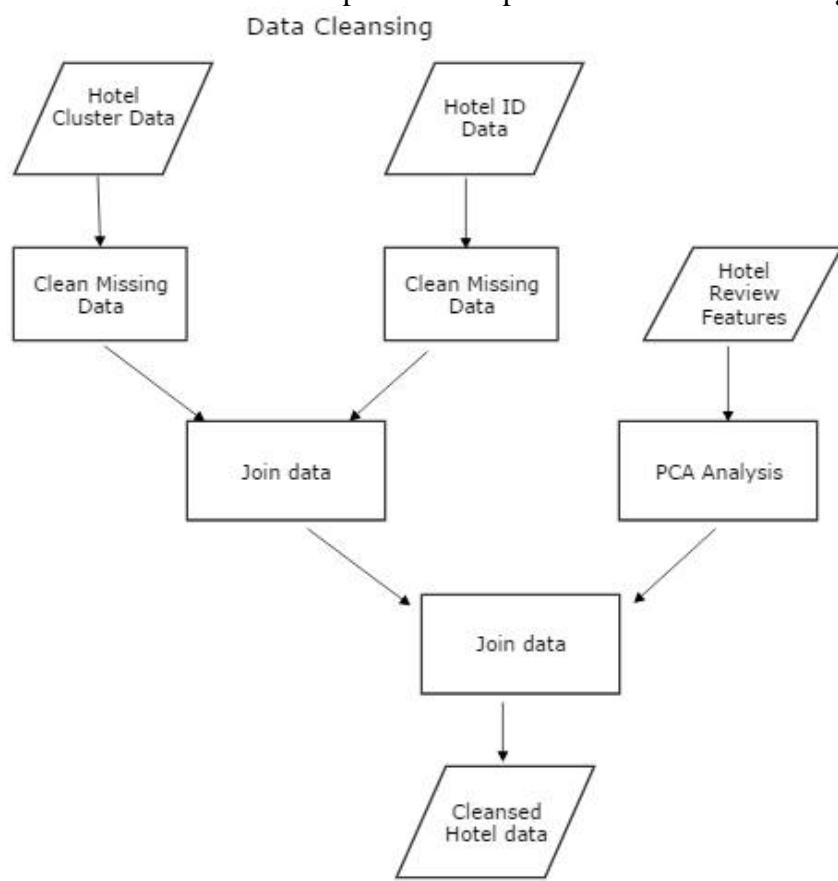
Graph4 – Visitor history star rating per hotel id



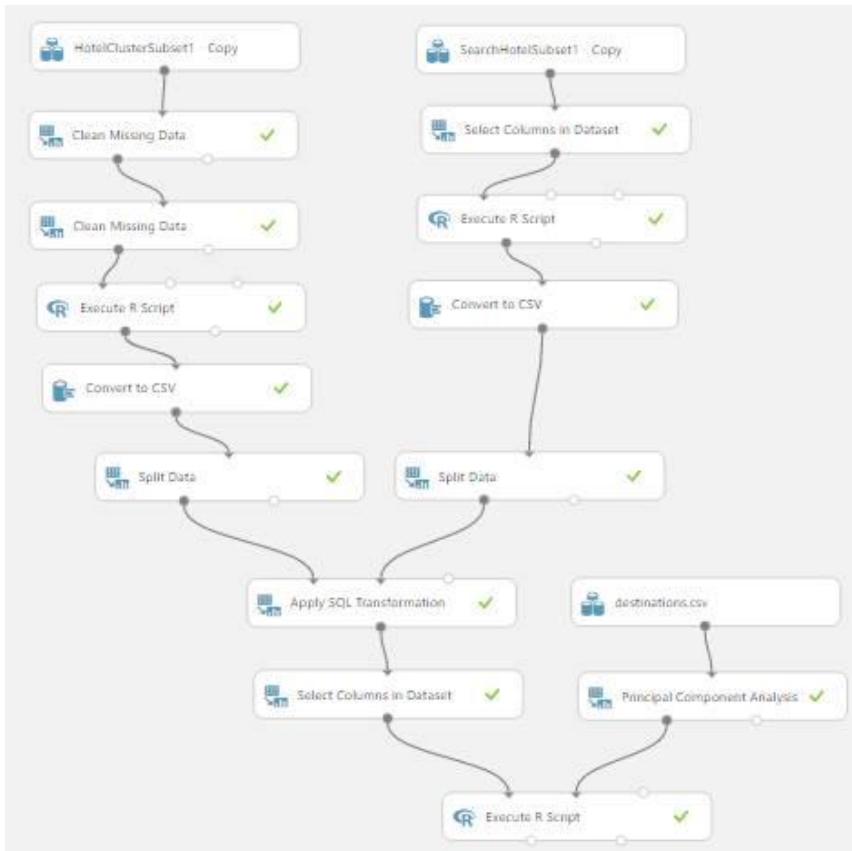
Graph5 – Tableau story**Hotel Cluster and User Story**

1.3 Data Pre-Processing

Hotel cluster and hotel Id are present in separate files which are merged together based on destination Id.



Azure Data cleansing model:



1.3.1 Missing Data Treatment

- In HotelCluster dataset approximately 0.12% data is missing for checkin and checkout data. So remove the rows with missing data.
- Column Orgi_destination_distance has 35% data as missing rows. One of the reasons it could be because the distance could not be calculated which can be because of presence of oceans in the origination and destination. Hence we decided to replace the missing values with the largest possible value present in the dataset.
- In visitor_hist_starrating and visitor_hist_adr_usd we have replaced all the NULL with '0' and ignored it later in our models.
- In search Hotel Dataset, column gross_booking_usd has rows with NULL value. To treat it calculate the formula for gross booking, which we found out as:

$$(\text{Price_usd} * \text{No. of rooms} * \text{length of stay}) + 15\% \text{ tax}$$

```

  · #replace orig_destination_distance null values with largest distance present in the column
  · maximumDistance <- max(searchHotelData$orig_destination_distance)
  · searchHotelData$orig_destination_distance <- ifelse((searchHotelData$orig_destination_distance=='NULL'),maximumDistance,searchHotelData$orig_destination_distance)

  · searchHotelData$gross_bookings_usd <- ifelse((searchHotelData$gross_bookings_usd=='NULL'),(
    · ((searchHotelData$price_usd * searchHotelData$srch_room_count) + (searchHotelData$price_usd * 0.15)
    ·   * searchHotelData$srch_length_of_stay),searchHotelData$gross_bookings_usd)
  ·

```

- Extract date, year, month, day and hour values for each date
- Create a column for Day of the week by assuming sun to Sat is represented in 0-6 format
- Create a column for IsWeekday by assuming 1-yes, 0-No
- Create a column for Booking_window by calculating the time difference between checkin date and booking date
- Create a column for stay_length by calculating the time difference between checkout date and checkin date

R Script

```

6 HotelClusterData$year <- as.numeric(format(as.Date(HotelClusterData$date_time,format = '%m/%d/%Y'), format= "%Y"))
7 HotelClusterData$month <- as.numeric(format(as.Date(HotelClusterData$date_time,format = '%m/%d/%Y'), format= "%m"))
8 HotelClusterData$day <- as.numeric(format(as.Date(HotelClusterData$date_time,format = '%m/%d/%Y'), format= "%d"))
9 HotelClusterData$hour <- format(as.POSIXct(HotelClusterData$date_time),'%H')
10 #monthlyNetworkData$hour <- as.numeric(format(as.Date(monthlyNetworkData$time,format = '%H:%M:%S'), format= "%H"))
11 |
12 #Get Day of the week--- sun to Sat in 0-6 format
13 require(lubridate)
14 HotelClusterData$DayOfWeek = wday(as.Date(HotelClusterData$date_time,format = '%Y-%m-%d'))-1
15
16 #IsWeekday ----- 1-yes, 0-No
17 require(timeDate)
18 HotelClusterData$Weekday <- ifelse(isWeekday(as.Date(HotelClusterData$date_time,format = '%Y-%m-%d')), wday = 1:5), 1, 0)
19
20 #HotelClusterData$Booking_Window <- difftime(HotelClusterData$srch_ci, HotelClusterData$date_time, units = "days") # days
21
22 HotelClusterData$Booking_Window <- as.Date(as.character(HotelClusterData$srch_ci), format = "%Y-%m-%d")-
23   as.Date(as.character(HotelClusterData$date_time), format = "%Y-%m-%d")
24 HotelClusterData$Stay_Length <- as.Date(as.character(HotelClusterData$srch_co), format = "%Y-%m-%d")-
25   as.Date(as.character(HotelClusterData$srch_ci), format = "%Y-%m-%d")

```

- Join the two dataset on the basis of destination ID.

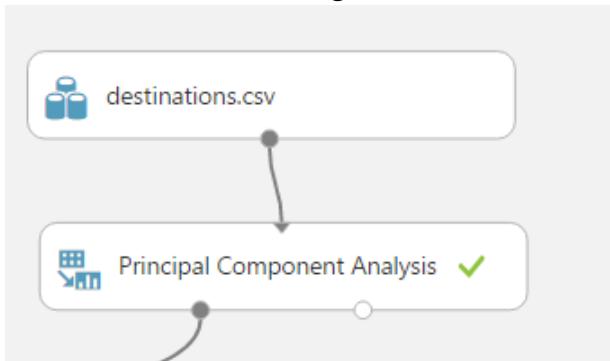
```
select * from t1 join t2 on t1.srch_destination_id = t2.srch_destination_id;
```

1.3.2 Feature Selection

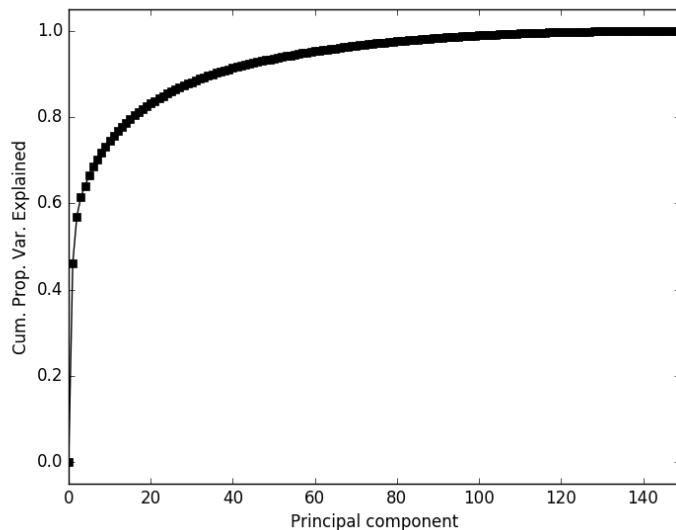
- Each of the features of the destination data set is actually a column from “d1” to “d149,” which means we don’t really know what the data within those columns mean. Furthermore, adding 149 columns of data to our training set for prediction will not be

particularly useful and will, more than likely, result in overfitting if we try to use them raw.

- One way to deal with that is to identify a small number of components—say, 5 of them—that explain a large amount of the variance among the 149 features.
- Apply **Principal component analysis (PCA)** on the destination file which contains the data on hotel rating features.

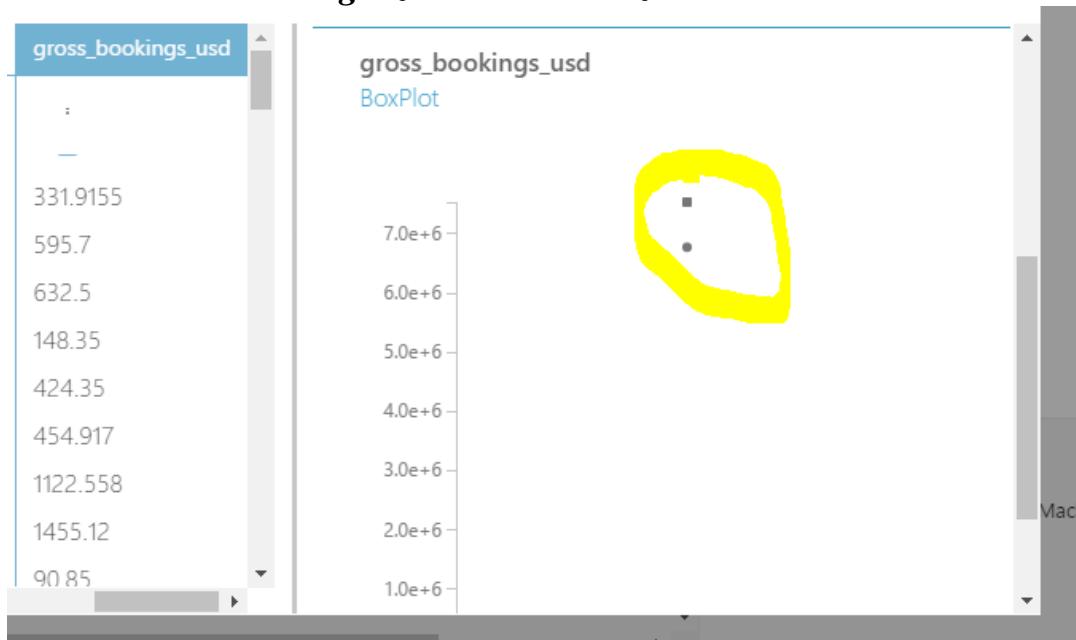
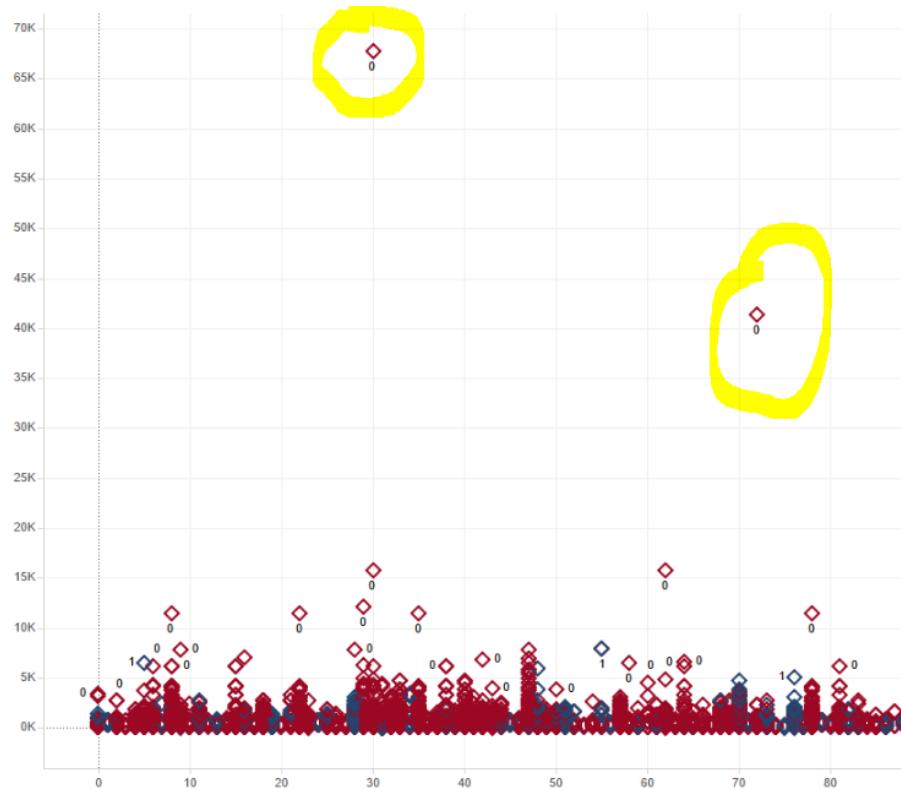


- Join the PCA analysed dataset to the hotel dataset.



1.3.3 Outliers detection and removal

- An **outlier** is an observation that appears to deviate markedly from other observations in the sample. Identification of potential **outliers** is important.
- Perform **outlier detection** on the combined dataset using Azure boxplot and tableau.

Outlier detection using Azure ML visualization:***Outlier detection using Tableau visualization:***

- Remove the outliers using execute R Script in Azure.

R Script

```

14 outlier.dataset <- c();
15 for(i in 1:nrow(hotelData)) {
16   if(((hotelData$gross_bookings_usd[i]) < (mean(hotelData$gross_bookings_usd) - (1.5)*sd(hotelData$gross_bookings_usd)) | 
17     (hotelData$gross_bookings_usd[i]) > (mean(hotelData$gross_bookings_usd) + (1.5)*sd(hotelData$gross_bookings_usd))) &
18     ((hotelData$srch_adults_cnt[i]) < (mean(hotelData$srch_adults_cnt) - (1.5)*sd(hotelData$srch_adults_cnt)) | 
19     (hotelData$srch_adults_cnt[i]) > (mean(hotelData$srch_adults_cnt) + (1.5)*sd(hotelData$srch_adults_cnt))) ) {
20     #&
21     #((hotelData$srch_children_cnt[i]) < (mean(hotelData$srch_children_cnt) - (1.5)*sd(hotelData$srch_children_cnt)) | 
22     # (hotelData$srch_children_cnt[i]) > (mean(hotelData$srch_children_cnt) + (1.5)*sd(hotelData$srch_children_cnt)) &
23     #((hotelData$srch_rm_cnt[i]) < (mean(hotelData$srch_rm_cnt) - (1.5)*sd(hotelData$srch_rm_cnt)) | 
24     #(hotelData$srch_rm_cnt[i]) > (mean(hotelData$srch_rm_cnt) + (1.5)*sd(hotelData$srch_rm_cnt))) &
25     #((hotelData$prop_country_id[i]) < (mean(hotelData$prop_country_id) - (1.5)*sd(hotelData$prop_country_id)) | 
26     #(hotelData$prop_country_id[i]) > (mean(hotelData$prop_country_id) + (1.5)*sd(hotelData$prop_country_id))) {
27     count = count + 1;
28     outlier.dataset <- c(outlier.dataset, i)
29   }
30 }
31 #outliers detected - 11
32 #removing the outlier from cleaned dataset
33 hotelDataNew <- data.frame(hotelData[-outlier.dataset,])

```

Azure Modules:

- **Split Data** - To split a dataset into two equal parts, just add the Split Data module after the dataset without no other changes. By default, the module splits the dataset in two equal parts. For data with an odd number of rows, the second output gets the remainder.

Split Data

Splitting mode

Split Rows

Fraction of rows in the first o...

0.75

 Randomized split

Random seed

0

Stratified split

False

START TIME 8/5/2016 3:53:...

END TIME 8/5/2016 3:53:...

ELAPSED TIME 0:00:00.000

- **Train Model** - Training a classification or regression model is a kind of *supervised machine learning*. That means you must provide a dataset that contains historical data from which to learn patterns. The data should contain both the outcome you are trying to predict,

and related factors (variables). The machine learning model uses the data to extract statistical patterns and build a model.

When you configure Train Model, you must also connect an already configured model, such as a regression algorithm, decision tree model, or other machine learning module.

- **Score Model** - Score Model is used to generate predictions using a trained classification or regression model. The predicted value can be in many different formats, depending on the model and your input data: If you are using a classification model to create the scores, Score Model outputs a predicted value for the class, as well as the probability of the predicted value. For regression models, Score Model generates just the predicted numeric value.
- **Evaluate Model** - Evaluate Model is used to measure the accuracy of a trained classification model or regression model. You provide a dataset containing scores generated from a trained model, and the Evaluate Model module computes a set of industry-standard evaluation metrics. The metrics returned by Evaluate Model depend on the type of model that you are evaluating

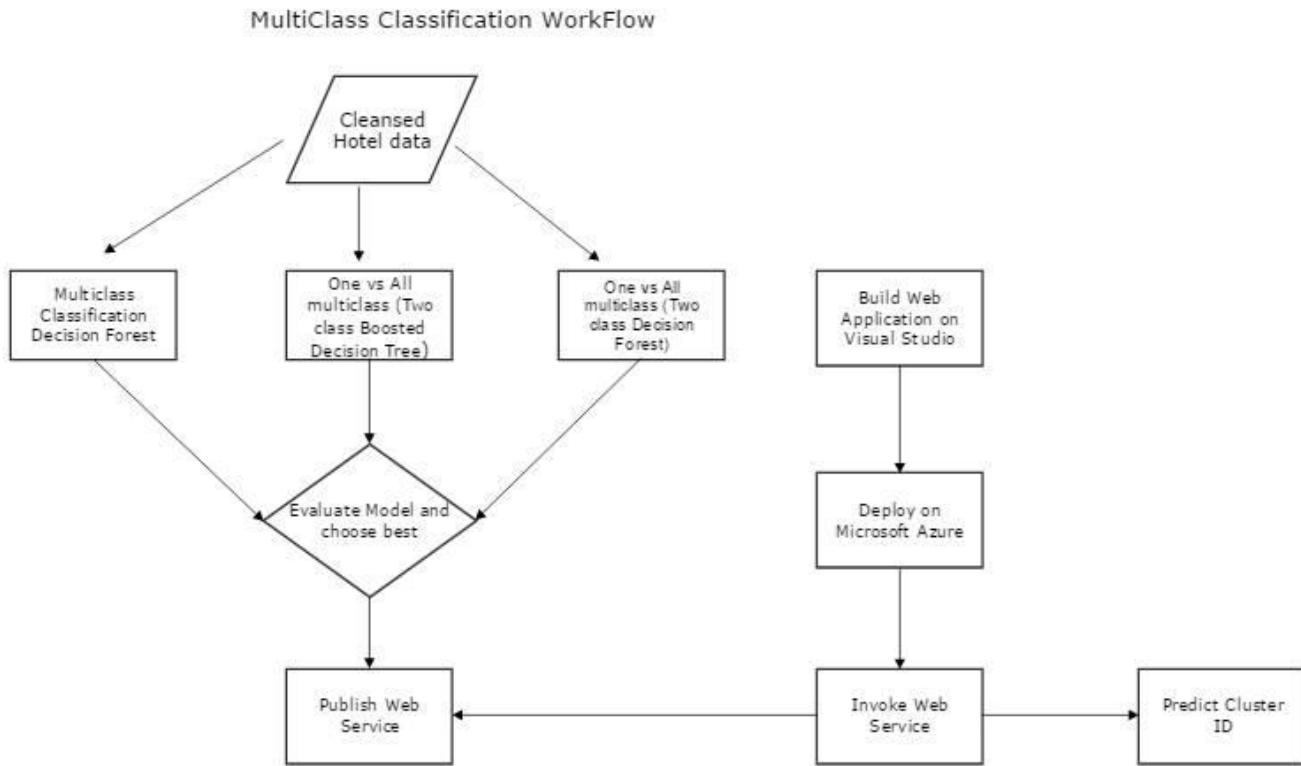
1.4 Classification Models

We build classification model to predict the booking outcome (hotel cluster) for a user event, based on their search and other attributes associated with that user event.

1.4.1 Overall Design

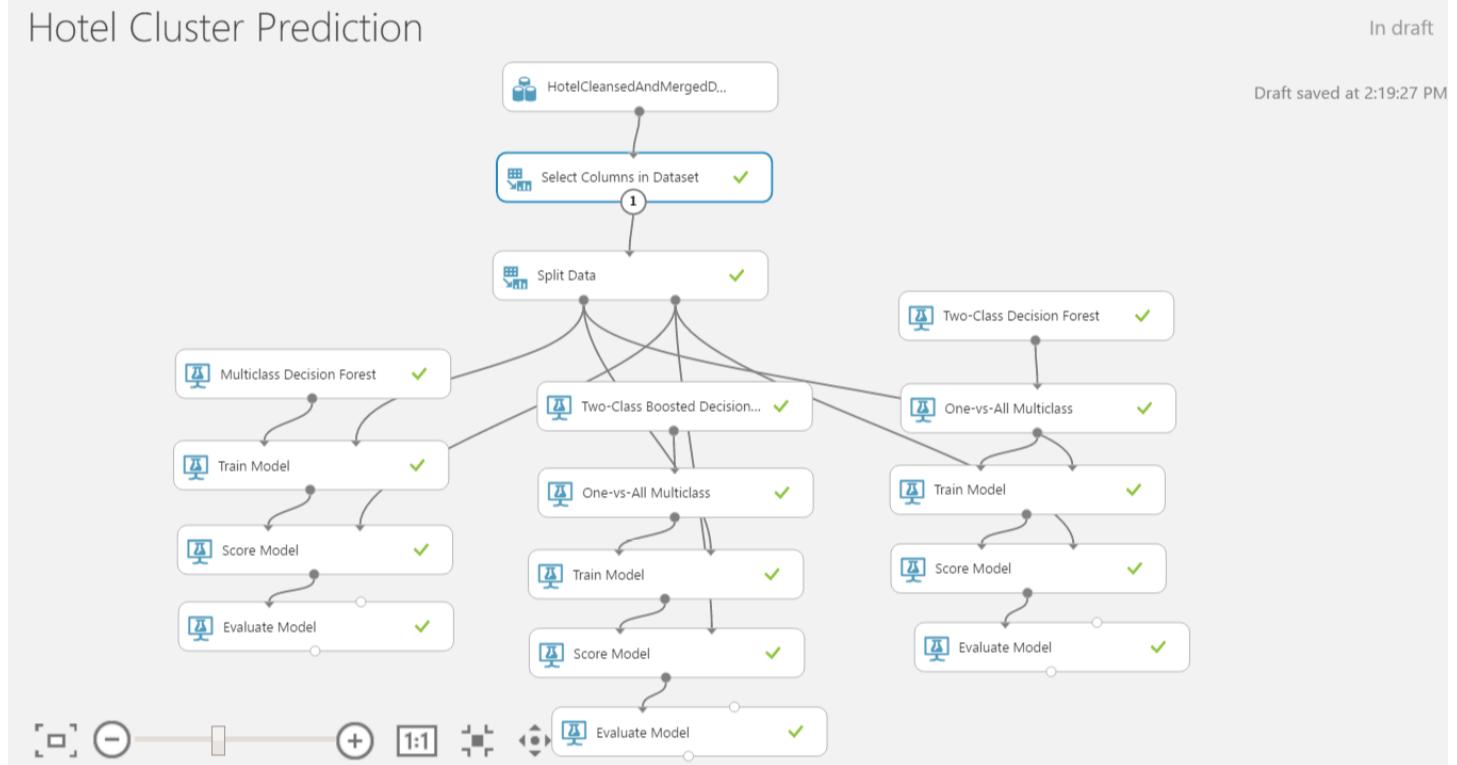
- ❖ Read the hotel recommendation cleansed and merged data.
- ❖ Select only the columns that have high coefficient value and exclude the rest.
- ❖ Split the data into train and test by 75% and 25% respectively.
- ❖ Implement Multi class Decision Forest, Two class Decision Forest and One-vs- All Multiclass, Two class Boosted Decision Tree and One-vs- All Multiclass.
- ❖ Parameters used in each model is elaborated in the further module below.
- ❖ Compare the models through the overall accuracy, average accuracy, macro-averaged precision, macro-averaged recall and choose the best model.
- ❖ Publish the best classification model created for this data as a web service.
- ❖ Build a web application using visual studio and deploy it on Microsoft Azure.
- ❖ This web application is used as an interface to invoke the web service expose through Azure.

- ❖ Predictions are performed through this web application to the booking outcome (hotel cluster) for a user event, based on their search and other attributes associated with that user event.



1.4.2 Azure models

Below is the hotel recommendation system classification model created for predicting the hotel cluster based on user event.



1.4.3 Multiclass Decision Forest

A machine learning model based on the decision forest algorithm. The decision forest algorithm is an ensemble learning method for classification. The algorithm works by building multiple decision trees and then voting on the most popular output class. Voting is a form of aggregation, in which each tree in a classification decision forest outputs a non-normalized frequency histogram of labels. The aggregation process sums these histograms and normalizes the result to get the “probabilities” for each label. The trees that have high prediction confidence will have a greater weight in the final decision of the ensemble. Decision trees in general are non-parametric models, meaning they support data with varied distributions. In each tree, a sequence of simple tests is run for each class, increasing the levels of a tree structure until a leaf node (decision) is reached. Decision trees have many advantages: They can represent non-linear decision boundaries. They are efficient in computation and memory usage during training and prediction. They perform integrated feature selection and classification. They are resilient in the presence of noisy features. The decision forest classifier in Azure Machine Learning Studio consists of an ensemble of decision trees. Generally, ensemble models provide better coverage and accuracy than single decision trees.

▲ Multiclass Decision Forest

Resampling method 



Create trainer mode 



Number of decision trees 



Maximum depth of the decisi... 



Number of random splits per... 



Minimum number of sample... 



Allow unknown values fo... 

Confusion Matrix and other performance metrics :

The following statistics are shown for our model:

Precision: Given all the predicted labels (for a given class X), how many instances were correctly predicted?

Recall: For all instances that should have a label X, how many of these were correctly captured?

The generalization to multi-class problems is to sum over rows / columns of the confusion matrix. Given that the matrix is oriented as above, i.e., that a given row of the matrix corresponds to specific value for the "truth", we have:

$$\text{Precision } i = \frac{M_{ii}}{\sum_j M_{ji}}$$

$$\text{Recall } i = \frac{M_{ii}}{\sum_j M_{ij}}$$

That is, precision is the fraction of events where we *correctly* declared i out of all instances where the algorithm declared i . Conversely, recall is the fraction of events where we correctly declared i out of all of the cases where the true state of the world is i

Hotel Cluster Prediction ➤ Evaluate Model ➤ Evaluation results

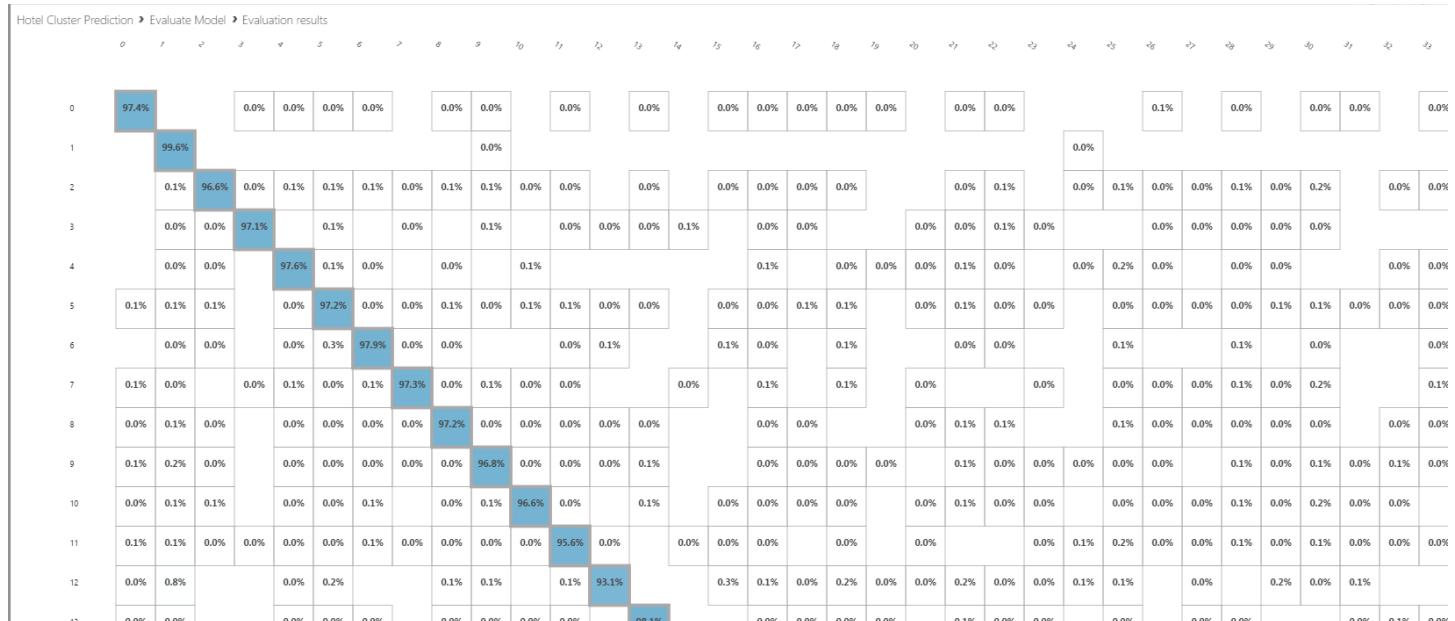
◀ Metrics

Overall accuracy	0.970708
Average accuracy	0.999414
Micro-averaged precision	0.970708
Macro-averaged precision	0.966281
Micro-averaged recall	0.970708
Macro-averaged recall	0.966178

The Need for a Confusion Matrix

Apart from helping with computing precision and recall, it is always important to look at the confusion matrix to analyze your results as it also gives you very strong clues as to where your classifier is going wrong.

Confusion Matrix



1.4.4 One – vs – All Multiclass

Azure ML Studio also provides a module called One-vs-All Multiclass which can use any binary classifier as an input to solve a multi-class classification problem, based on this [one-vs-](#)

all method. Therefore, as the second model for comparison, we used a binary classification module, **Two-Class Boosted Decision Tree**, and connected it to the **One-vs-All Multiclass** module. Third model for comparison, we used a binary classification module, **Two-Class Decision Forest**, and connected it to the **One-vs-All Multiclass** module.

1.4.4.1 Two class Boosted Decision Tree

A boosted decision tree is an ensemble learning method in which the second tree corrects for the errors of the first tree, the third tree corrects for the errors of the first and second trees, and so forth. Predictions are based on the entire ensemble of trees together that makes the prediction. Generally, when properly configured, boosted decision trees are the easiest methods with which to get top performance on a wide variety of machine learning tasks. However, they are also one of the more memory-intensive learners, and the current implementation holds everything in memory; therefore, a boosted decision tree model might not be able to process the very large datasets that some linear learners can handle.

▲ Two-Class Boosted Decision Tree

Create trainer mode

Single Parameter ▾

Maximum number of leaves ...

Minimum number of sample...

Learning rate

Number of trees constructed

Random number seed

Confusion Matrix and other performance metrics :

The following statistics are shown for our model:

Precision: Given all the predicted labels (for a given class X), how many instances were correctly predicted?

Recall: For all instances that should have a label X, how many of these were correctly captured?

The generalization to multi-class problems is to sum over rows / columns of the confusion matrix. Given that the matrix is oriented as above, i.e., that a given row of the matrix corresponds to specific value for the "truth", we have:

Precision $i = M_{ii} / \sum_j M_{ji}$

$$\text{Recall } i = M_{ii} / \sum_j M_{ij}$$

That is, precision is the fraction of events where we *correctly* declared i_i out of all instances where the algorithm declared i_i . Conversely, recall is the fraction of events where we correctly declared i_i out of all of the cases where the true state of the world is i .

Hotel Cluster Prediction ➤ Evaluate Model ➤ Evaluation results

► Metrics

Overall accuracy	0.7463
Average accuracy	0.994926
Micro-averaged precision	0.7463
Macro-averaged precision	0.787456
Micro-averaged recall	0.7463
Macro-averaged recall	0.731311

The Need for a Confusion Matrix

Apart from helping with computing precision and recall, it is always important to look at the confusion matrix to analyze your results as it also gives you very strong clues as to where your classifier is going wrong.

Confusion Matrix

1.4.4.2 Two class Decision Forest

The decision forest algorithm is an ensemble learning method for classification. The algorithm works by building multiple decision trees and then voting on the most popular output class. Voting is a form of aggregation, in which each tree in a classification decision forest outputs a non-normalized frequency histogram of labels. The aggregation process sums these histograms and normalizes the result to get the “probabilities” for each label. The trees that have high prediction confidence will have a greater weight in the final decision of the ensemble.

▲ Two-Class Decision Forest

Resampling method 

Bagging 

Create trainer mode 

Single Parameter 

Number of decision tre... 

8 

Maximum depth of the... 

32 

Number of random spl... 

128 

Minimum number of s... 

1 

Allow unknown val... 

Confusion Matrix and other performance metrics :

The following statistics are shown for our model:

Precision: Given all the predicted labels (for a given class X), how many instances were correctly predicted?

Recall: For all instances that should have a label X, how many of these were correctly captured?

The generalization to multi-class problems is to sum over rows / columns of the confusion matrix. Given that the matrix is oriented as above, i.e., that a given row of the matrix corresponds to specific value for the "truth", we have:

$$\text{Precision } i = \frac{M_{ii}}{\sum_j M_{ij}}$$

$$\text{Recall } i = \frac{M_{ii}}{\sum_j M_{ji}}$$

That is, precision is the fraction of events where we *correctly* declared i out of all instances where the algorithm declared i . Conversely, recall is the fraction of events where we correctly declared i out of all of the cases where the true state of the world is i .

Hotel Cluster Prediction ➤ Evaluate Model ➤ Evaluation results

Metrics

Overall accuracy	0.914582
Average accuracy	0.998292
Micro-averaged precision	0.914582
Macro-averaged precision	0.906533
Micro-averaged recall	0.914582
Macro-averaged recall	0.902382

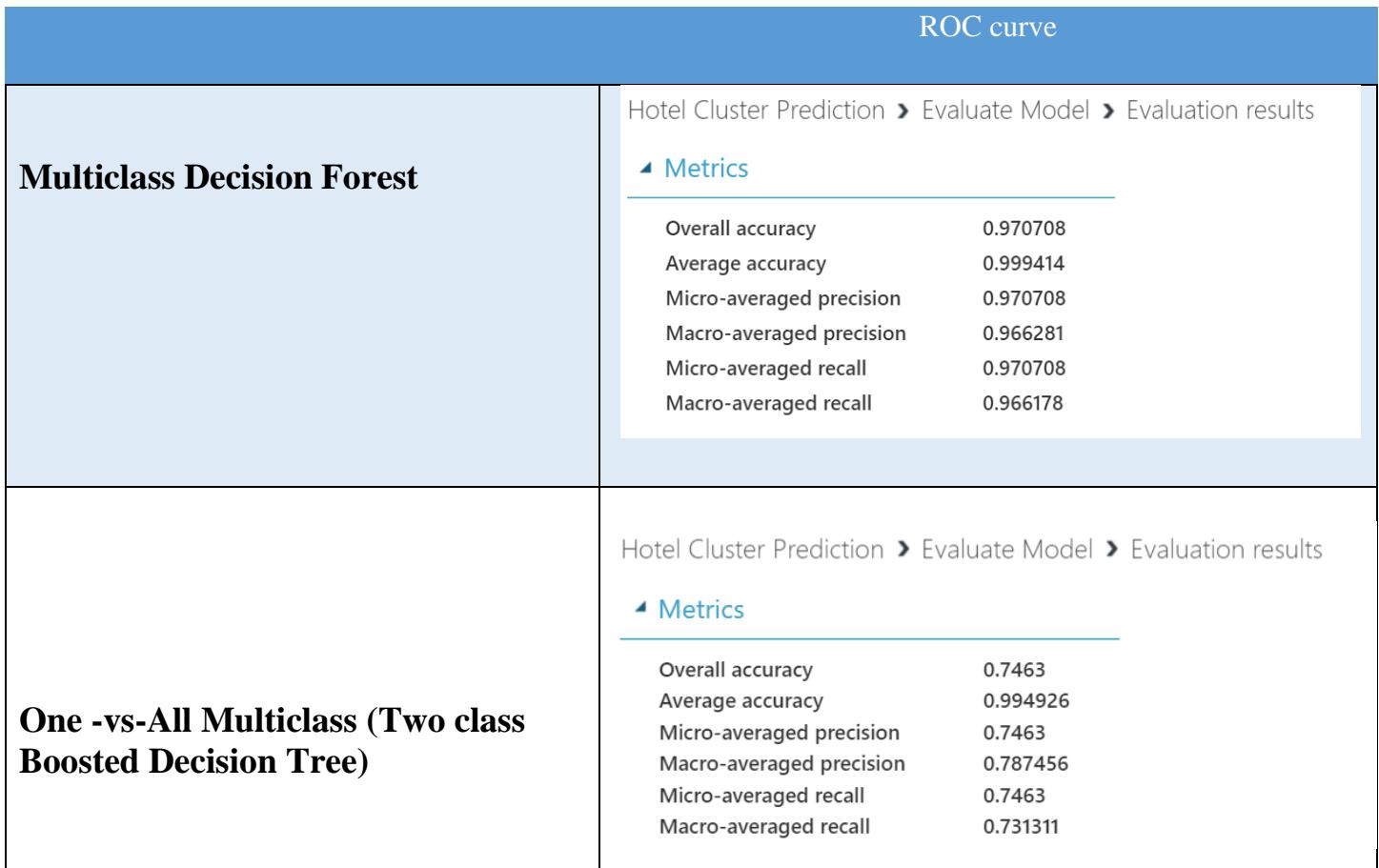
The Need for a Confusion Matrix

Apart from helping with computing precision and recall, it is always important to look at the confusion matrix to analyze your results as it also gives you very strong clues as to where your classifier is going wrong.

Confusion Matrix

Hotel Cluster Prediction > Evaluate Model > Evaluation results

1.4.5 Classification model comparison



One -vs-All Multiclass (Two class Decision Forest)	<p>Hotel Cluster Prediction ➔ Evaluate Model ➔ Evaluation results</p> <p>◀ Metrics</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <tbody> <tr> <td style="width: 40%;">Overall accuracy</td><td style="width: 60%;">0.914582</td></tr> <tr> <td>Average accuracy</td><td>0.998292</td></tr> <tr> <td>Micro-averaged precision</td><td>0.914582</td></tr> <tr> <td>Macro-averaged precision</td><td>0.906533</td></tr> <tr> <td>Micro-averaged recall</td><td>0.914582</td></tr> <tr> <td>Macro-averaged recall</td><td>0.902382</td></tr> </tbody> </table>	Overall accuracy	0.914582	Average accuracy	0.998292	Micro-averaged precision	0.914582	Macro-averaged precision	0.906533	Micro-averaged recall	0.914582	Macro-averaged recall	0.902382
Overall accuracy	0.914582												
Average accuracy	0.998292												
Micro-averaged precision	0.914582												
Macro-averaged precision	0.906533												
Micro-averaged recall	0.914582												
Macro-averaged recall	0.902382												

Conclusion:

In evaluating multi-class classification problems, we often think that the only way to evaluate performance is by computing the accuracy which is the proportion or percentage of correctly predicted labels over all predictions.

However, we can always compute precision and recall for each class label and analyze the individual performance on class labels or average the values to get the overall precision and recall.

Accuracy alone is sometimes quite misleading as you may have a model with relatively 'high' accuracy with the model predicting the '***not so important***' class labels fairly accurately (e.g. "unknown bucket") but the model may be making all sorts of mistakes on the classes that are actually critical to the application.

Precision: Given all the predicted labels (for a given class X), how many instances were correctly predicted?

Recall: For all instances that should have a label X, how many of these were correctly captured?

We choose Multiclass Decision Forest, which means that for precision, out of the times label '1' was predicted, 99.6% of the time the system was in fact correct. And for recall, it means that out of all the times label '1' should have been predicted only 99.6% of the labels were correctly predicted. (Please refer confusion matrix)

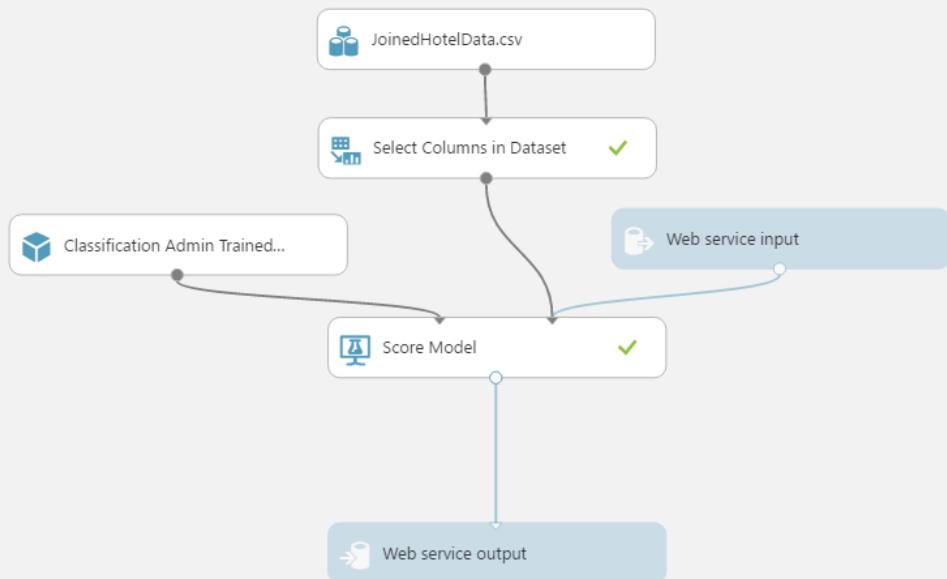
In essence, the more zeroes or smaller the numbers on all cells in the confusion but the diagonal, the better your classifier is doing.

Hence, Multiclass Decision Forest has better accuracy, precision and recall rate than all other model.

1.4.6 Web Service

- Once the classification model is ready, we set up **Web Service**.
- The model we trained is saved as a single **Trained Model** module into the module palette to the left of the experiment canvas (you can find it under **Trained Models**).
- Modules that were used for training are removed. Specifically:
 - Multi class decision forest
 - Train Model
 - Split Data
- Then we added the saved trained model back into the experiment.
- **Web service input** and **Web service output** modules are added.

Hotel Classification Admin Webservice Deployed



- Now run the model and publish the web service.

hotel classification admin

DASHBOARD CONFIGURATION

General

Published experiment

[View snapshot](#) [View latest](#)

Description

No description provided for this web service.

API key

d/jnA+4Z0W0TRIBRhUaUIN6BeaRPZ09urfmX3/rWADjBajgBQfEa4mLTkJzAri0pvfr2ELMO4sSYgVHiuogvg=



Default Endpoint

[API HELP PAGE](#)

TEST

APPS

LAST UPDATED



[REQUEST/RESPONSE](#)

Test

Excel 2013 or later | Excel 2010 or earlier workbook

8/19/2016 12:03:37 AM

[BATCH EXECUTION](#)

Excel 2013 or later workbook

8/19/2016 12:03:37 AM

- On running the web service, we get the following form which can be used to invoke the web service and do prediction.

Test Hotel Classification Admin Service

Enter data to predict

USER_LOCATION_CITY



USER_ID

IS_MOBILE

IS_PACKAGE

SRCH_ADULTS_CNT



- ❖ Web application making call to the API and invoking the cluster web service.

The screenshot shows the Microsoft Visual Studio interface. The title bar reads "View all photos - Microsoft Visual Studio". The menu bar includes File, Edit, View, Project, Build, Debug, Team, Tools, Test, Analyze, Window, and Help. The toolbar has icons for Quick Launch, View, Project, Build, Debug, Team, Tools, Test, Analyze, Window, and Help. The status bar at the bottom shows "Dheyansh Srivastava" and "100 %".

The code editor displays a C# file named "MyClusters.cs" with the following content:

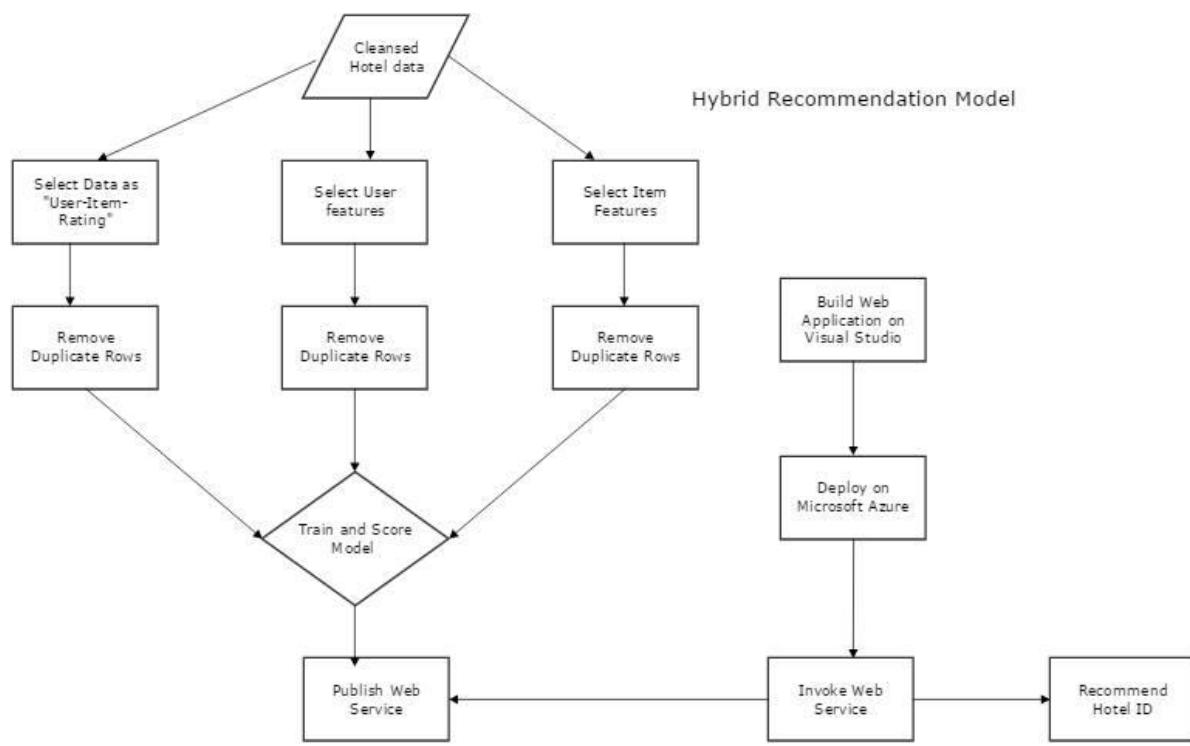
```
model.Stay_Length = (model.checkoutDate - model.checkinDate).TotalDays + ".00:00:00";

using (var client = new HttpClient())
{
    var scoreRequest = new
    {
        Inputs = new Dictionary<string, StringTable>()
    {
        "input1",
        new StringTable()
    {
        ColumnNames = new string[] {"user_location_city", "user_id", "is_mobile", "is_pac"},
        Values = new string[,] { { model.user_location_city.ToString(), model.user_id.ToString() } }
    }
},
        GlobalParameters = new Dictionary<string, string>()
    {
    }
};
const string apiKey = "d/jnA+4Z0W0TRIBRhUaUIN6BearPZ09urfMx3/rWADjBajgBQfEa4mLTKJzAri0pvfr2ELMO4";
client.DefaultRequestHeaders.Authorization = new AuthenticationHeaderValue("Bearer", apiKey);
client.BaseAddress = new Uri("https://ussouthcentral.services.azureml.net/workspaces/b78cb51f8e19");
// WARNING: The 'await' statement below can result in a deadlock if you are calling this code from a UI thread.
```

The Solution Explorer on the right lists the project structure:

- Images
- bootstrap.css
- bootstrap.min.css
- + HotelIDAndReviews.csv
- + rangeslider.css
- + rangeslider.js
- + rangeslider.min.js
- + Report.pdf
- + SingleUser.csv
- Site.css
- style.css
- UseridHotel.csv
- Controllers
- + AccountController.cs
- + AdminController.cs
- + AzureBlobDataReference.cs
- + ExecuteRcs
- + Helper.cs
- + HomeController.cs
- + ManageController.cs
- + MyClusters.cs
- + MyData.cs
- + RecommendationBatch.cs
- fonts

1.5 Recommendation System



1.5.1 Overall Design

- ❖ This section explains the implementation and the application of a Recommendation system built around the dataset.
- ❖ The main aim of a recommendation system is to recommend one or more items to users of the system.
- ❖ There are two approaches to recommender systems, content-based and collaborative filtering. To implement a recommender system, we used the Matchbox Recommender on ML Studio.
- ❖ The Matchbox recommender on ML Studio combines both the approaches, Collaborative filtering and Content based thus also called as a hybrid recommender.

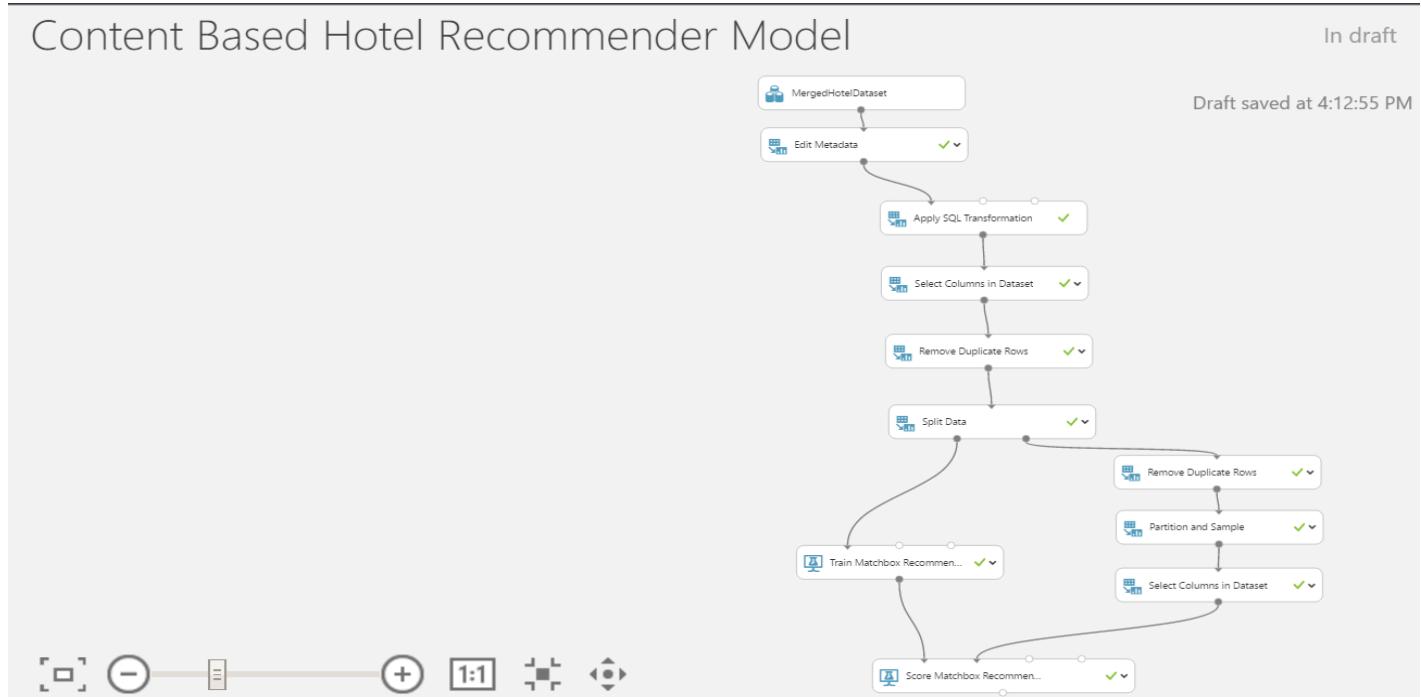
Recommender systems are classified according to the technique used to create the recommendation (fill the blanks in the utility matrix):

- ⊕ Content-based systems examine properties of the items recommended and offer similar items
- ⊕ Collaborative filtering (CF) systems recommend items based on similarity measures between users and/or items. The items recommended to a user are those preferred by similar users
- ⊕ Hybrid mixing both previous approaches

1.5.2 Content Based Recommendation System

1.5.2.1 Build model

- ❖ We only want three columns, since the MatchBox recommender that we will be using for training a recommender only takes in a dataset of triples: (user, item, rating).
- ❖ Translating it into our case, we need a dataset of the following three columns: UserId, PropId and Rating.
- ❖ The matchbox recommender has a limit of 0-99 or 1-100 distribution of ratings so the model would not train.
- ❖ Thus we rounded up the ratings and thus now the range is from 1-5 in a R Script in MLStudio.



Training, Scoring and Evaluating Matchbox Recommender:

- ❖ To train the matchbox recommender, we split the data into test and train datasets using the recommender split option in the Split data module.
- ❖ Split the data in the proportion of 75% and 25%.
- ❖ Then we Add the Train MatchBox Recommender module and connect it the training data which returns a trained Matchbox recommender. We then used the Score Matchbox Recommender module, which creates recommendations for the different users.

Content Based Hotel Recommender We... ➤ Score Matchbox Recommender

rows columns
100 4

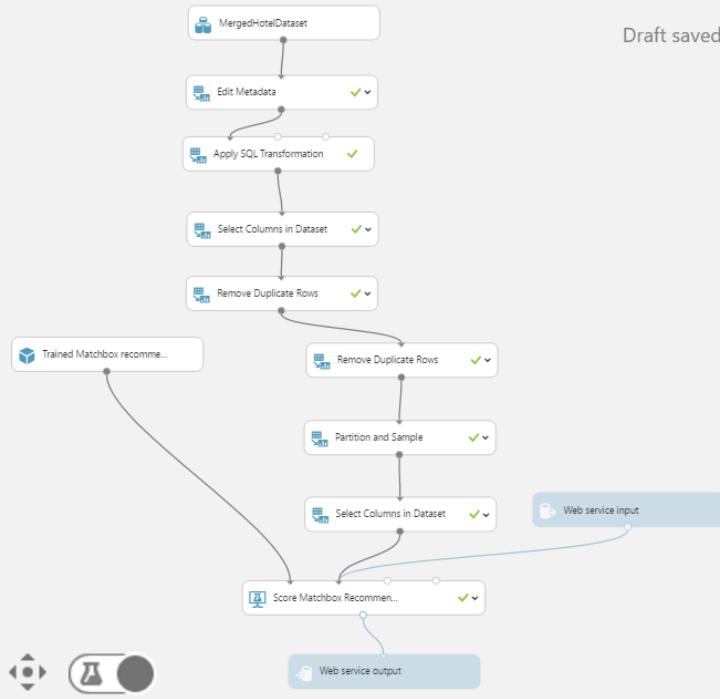
	User	Item 1	Item 2	Item 3
view as				
6	118321	1506	60644	
11	118321	1506	60644	
12	118321	1506	60644	
13	118321	1506	60644	
15	118321	1506	60644	
33	118321	91863	60644	
40	118321	1506	60644	
67	118321	1506	60644	

1.5.2.2 Publish Content-Based Filtering Model as a Web Service

- ❖ The model we trained is saved as a single **Trained Model** module into the module palette to the left of the experiment canvas (you can find it under **Trained Models**).
- ❖ Modules that were used for training are removed. Specifically:
 - Train Recommender model
 - Split Data
- ❖ Then we added the saved trained model back into the experiment.
- ❖ **Web service input** and **Web service output** modules are added.

Content Based Hotel Recommender WebService Deployed

Finished running ✓



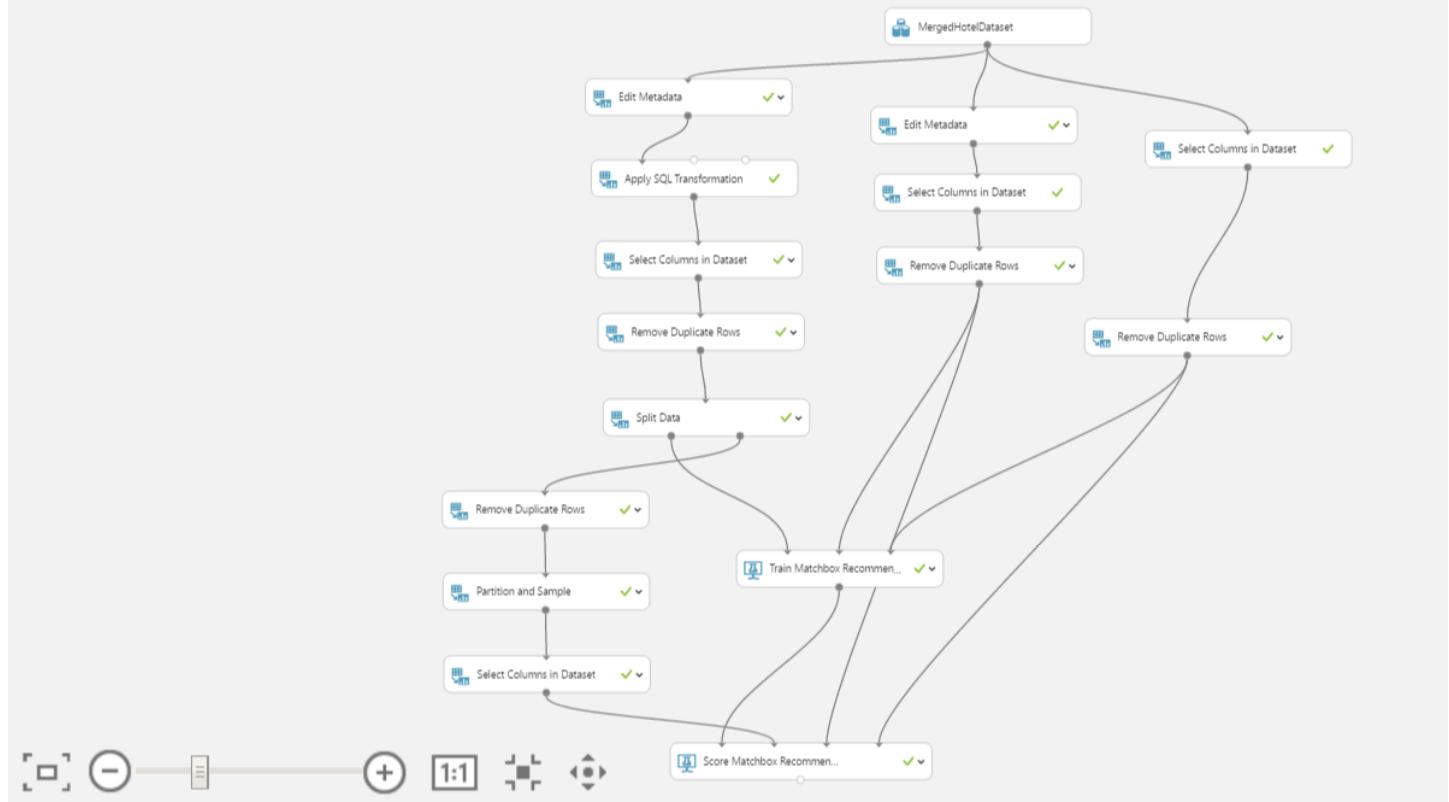
1.5.3 Hybrid Based Recommendation System

We now want to extend the content-/rating-based filtering approach to a hybrid recommender. This can be done by integrating user as well as item features - the remaining two inputs of the module Train Matchbox Recommender. User features encompass more information on the users, such as demographic information, while item features contain information on the hotels, e.g. usd, review score, destination Id

1.5.2.3 Build model

- ❖ We only want three columns, since the MatchBox recommender that we will be using for training a recommender only takes in a dataset of triples: (user, item, rating).
- ❖ Translating it into our case, we need a dataset of the following three columns: UserId, PropId and Rating.
- ❖ The matchbox recommender has a limit of 0-99 or 1-100 distribution of ratings so the model would not train.
- ❖ Thus we rounded up the ratings and thus now the range is from 1-5 in a R Script in MLStudio.

Hybrid Hotel Recommender Model



Training, Scoring and Evaluating Matchbox Recommender:

- ❖ To train the matchbox recommender, we split the data into test and train datasets using the recommender split option in the Split data module.
- ❖ Split the data in the proportion of 75% and 25%.
- ❖ Then we Add the Train MatchBox Recommender module and connect it the training data which returns a trained Matchbox recommender. We then used the Score Matchbox Recommender module, which creates recommendations for the different users.

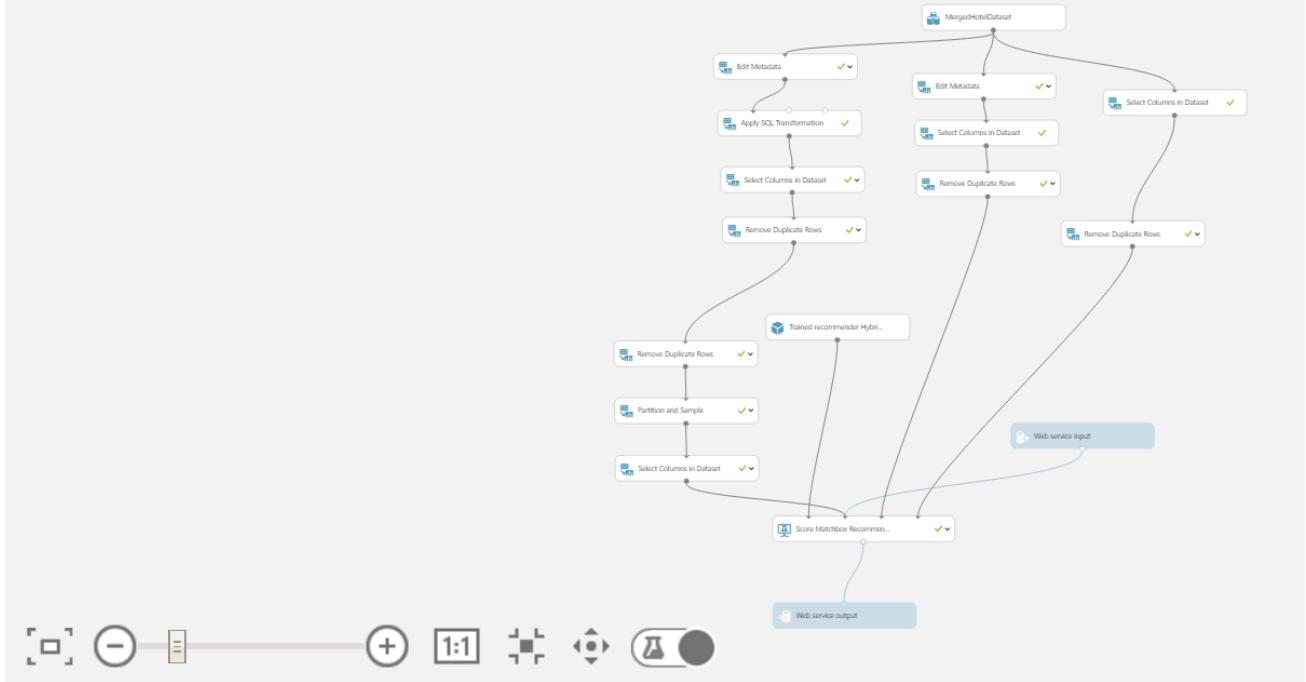
Hybrid Hotel Recommender WebService... ➤ Score Matchbox Recommender ➤

rows	columns	User	Item 1	Item 2	Item 3
100	4				
view as					
6	130585	125401	67657		
11	11905	67657	110744		
12	130585	125401	67657		
13	130585	125401	67657		
15	67657	130585	107820		
33	67657	130585	107820		
40	11905	67657	110744		
67	67657	130585	107820		

1.5.2.4 Publish Hybrid recommendation model as a Web Service

- ❖ The model we trained is saved as a single **Trained Model** module into the module palette to the left of the experiment canvas (you can find it under **Trained Models**).
- ❖ Modules that were used for training are removed. Specifically:
 - Train Recommender model
 - Split Data
- ❖ Then we added the saved trained model back into the experiment.
- ❖ **Web service input** and **Web service output** modules are added.

Hybrid Hotel Recommender WebService Deployed



- ❖ Input from the user is read and is stored in Azure blob storage.

```

154     {
155         // How this works:
156         //
157         // 1. Assume the input is present in a local file (if the web service accepts input)
158         // 2. Upload the file to an Azure blob - you'd need an Azure storage account
159         // 3. Call the Batch Execution Service to process the data in the blob. Any output is written to
160         // 4. Download the output blob, if any, to local file
161         if (!Helper.writefile(userid))
162         {
163             return "writefileFailed";
164         }
165         const string BaseUrl = "https://ussouthcentral.services.azureml.net/workspaces/b78cb51f8e194159";
166
167         const string StorageAccountName = "adsfinalteam3"; // Replace this with your Azure Storage Account
168         const string StorageAccountKey = "p7nFH+ilX9emmXsVTwmSMGKVzBd4QcITNNSIPZCR1af+SyaGZHZEYg+KcyQ";
169         const string StorageContainerName = "adsfinaldata"; // Replace this with your Azure Storage Container
170         const string apiKey = "jUFXNd5oh4N@Vcw2Fjk1NfSgZqPgz+vCDUcQnNYQpBBHjf9eDPNMK8rQldlFwU5F1Psrgo";
171
172         // set a time out for polling status
173         const int TimeOutInMilliseconds = 600 * 1000; // Set a timeout of 2 minutes
174
175
176         string storageConnectionString = string.Format("DefaultEndpointsProtocol=https;AccountName={0};",
177
178             String pth = System.Web.HttpContext.Current.Server.MapPath("~/");
179             pth += "/Content/";
180             UploadFileToBlob(pth + "SingleUser" + userid + ".csv" /*Replace this with the location of your
181             "SingleUserBlob" + userid + ".csv" *//*Replace this with the name you would like to use for your
             */
        
```

- ❖ Batch execution service processes the data on blob and gives an output file which is again dumped on azure blob.
- ❖ Output is read from this file and displayed back to the user.

```

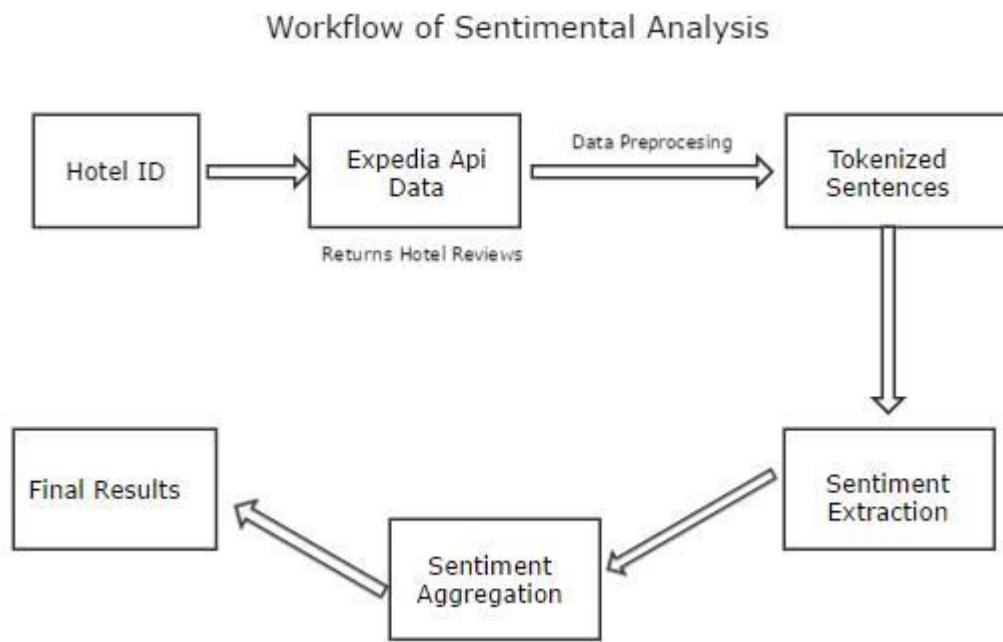
50         return View();
51     }
52 
53     [AllowAnonymous]
54     public ActionResult Result()
55     {
56 
57         ResultViewModel model = new ResultViewModel();
58         model.recommendations = new List<string>();
59         try
60         {
61             model.userid = User.Identity.GetUserId<int>();
62             ViewBag.Message = "Your Recommendations";
63             int userId = User.Identity.GetUserId<int>();
64             RecommendationBatch.InvokeBatchExecutionService(userId.ToString());
65             model.recommendations = Helper.readRecommendations(model.userid.ToString());
66         }
67         catch (Exception e) { }
68         return PartialView(model);
69     }
70 
71     [AllowAnonymous]
72     [HttpGet]
73     public ActionResult Result(int userid)
74     {
75         ResultViewModel model = new ResultViewModel();
76         if(userid==1)
77             model.userid = User.Identity.GetUserId<int>();
    }
  
```

1.6 Sentimental Analysis

Sentiment analysis (also known as **opinion mining**) refers to the use of natural language processing, text analysis and computational linguistics to identify and extract subjective information in source materials. Sentiment analysis is widely applied to reviews and social media for a variety of applications, ranging from marketing to customer service.

Generally speaking, sentiment analysis aims to determine the attitude of a speaker or a writer with respect to some topic or the overall contextual polarity of a document. The attitude may be his or her judgment or evaluation, affective state, or the intended emotional communication.

1.6.1 Overall Design



- ❖ We use a public API from Expedia to get reviews of all those hotels.

[Expedia](#) [Public APIs](#) [Partner APIs](#) [Sample Use Cases](#) [My Apps](#) [FAQs](#) [murali.a@husky.n...](#) [Logout](#)

Hotel Reviews

Description Try it out!

To get an API Key, Register or Login and go to My Apps. Adding APIs to your applications will authorize your key for use with those APIs.

Retrieve verified user reviews for a given hotel. All reviews are written by real Expedia customers who have stayed at the respective property.
Sample hotel ids are Hilton Seattle: 7910, Westin Seattle: 284304

Request Format

REST URL: `GET http://terminal2.expedia.com/x/reviews/hotels?hotelId={hotelid}&summary={issummarybool}&sortBy={sortstyle}&start={start}&items={items}&apikey={apikey}`

Example Requests

Return hotel reviews for the specified hotel, by hotel ID.
`http://terminal2.expedia.com/x/reviews/hotels?hotelId=234&apikey=[INSERT_KEY_HERE]`

Return a hotel review for the specified hotel, by review ID.
`http://terminal2.expedia.com/x/reviews/hotels?reviewId=18111063&apikey=[INSERT_KEY_HERE]`

 [Privacy Policy](#) | [Email Us](#)

- ❖ Code snippet to access the API and fetch all hotel reviews.

```

json_file <- paste("http://terminal2.expedia.com/x/reviews/hotels?hotelId=",hotelId,"&apiKey=JvMtGQVznEaKqvkEpYFZHSYEqINUwhH",sep="",collapse = "")

#Read data from file
json_data <- fromJSON(file=json_file)

reviewVector <- vector()
#iterate through the reviews of hotel and read the text
reviewObjectLength <- length(json_data$reviewDetails$reviewCollection$review)
for(reviewNumber in 1:reviewObjectLength)
{
  if(!length(json_data$reviewDetails$reviewCollection$review) == 0)
  {
    # print("step1 : no review")
    if(!is.null(json_data$reviewDetails$reviewCollection$review[[reviewNumber]]))
    {
      reviewVector <- c(reviewVector,json_data$reviewDetails$reviewCollection$review[[reviewNumber]][["reviewText"]])
    }
  }
}

```

- ❖ A word cloud helps us to visualize the most common words in the reviews and have a general feeling of the reviews.

```

if(!(length(reviewVector) == 0))
{
# print("step2")
#Converting vectors to matrix
hotelReviewMatrix <- matrix(reviewVector ,ncol = 1, byrow=TRUE)

# Create corpus
corpus=Corpus(VectorSource(reviewVector))

# Convert to lower-case
corpus=tm_map(corpus,tolower)

# Remove stopwords
corpus=tm_map(corpus,function(x) removeWords(x,stopwords()))

# convert corpus to a Plain Text Document
corpus=tm_map(corpus,PlainTextDocument)

col=brewer.pal(6,"Dark2")
wordcloud(corpus, min.freq=25, scale=c(5,2),rot.per = 0.25,
          random.color=T, max.word=45, random.order=F,colors=col)

```

 100%

Word cloud:

car hilton
airport nice
staff arrived.
stay

clean

- ❖ We will use lexicon based sentiment analysis. A list of positive and negative opinion words or sentiment words for English has been attached.

```
# print("step3")
#lexicon based sentiment analysis
positives= readLines("positive-words.txt")
negatives = readLines("negative-words.txt")

reviewScore <- data.frame()
for(row in 1:nrow(hotelReviewMatrix))
{
  review = hotelReviewMatrix[row,1]

  review = gsub("[[:punct:]]", "", review)      # remove punctuation
  review = gsub("[[:cntrl:]]", "", review)      # remove control characters
  #review = gsub('\d+', "", review)            # remove digits
  review = str_replace_all(review, "[^[:alnum:]]", " ")

  # Let's have error handling function when trying tolower
  tryToLower = function(x){
    # create missing value
    y = NA
    # tryCatch error
    try_error = tryCatch(tolower(x), error=function(e) e)
    # if not an error
    if (!inherits(try_error, "error"))
      y = tolower(x)
    # result
    return(y)
  }
}
```

- ❖ We calculate the score for each word in the review by matching it with positive and negative list of words. And finally an overall score is calculated which is more of a consolidated review of all the users of that hotel.

```
# use tryToLower with sapply
review = sapply(review, tryToLower)

require(stringr)
# split sentence into words with str_split function from stringr package
word_list = str_split(review, " ")
words = unlist(word_list)

# compare words to the dictionaries of positive & negative terms
positive.matches = match(words, positives)
negative.matches = match(words, negatives)

# get the position of the matched term or NA
# we just want a TRUE/FALSE
positive_matches = !is.na(positive.matches)
negative_matches = !is.na(negative.matches)

# final score
score = sum(positive_matches) - sum(negative_matches)
reviewScore = c(reviewScore, score)
```

1.6.2 Integration with Web Application

- ❖ This file with hotel Id and review score is stored on azure blob. Based on the score, a traveler rating is assigned to each hotel that the customer can view on log -in.

1.7 Web Application

We created a MVC ASP.net application on visual studio and then deployed the web application on Microsoft Azure.

Website link: <http://hoteladvisorsystem.azurewebsites.net/>

HotelSystem - Microsoft Visual Studio

File Edit View Project Build Debug Team Tools Test Analyze Window Help

Debug Any CPU Google Chrome

Server Explorer Toolbox Cloud Explorer

Result.cshtml Index.cshtml _Layout.cshtml Review.cshtml Classify.cshtml LoginPartial.cshtml Business.cshtml

```

9     }
10    int cnt = 0;
11  }
12  <link rel="stylesheet" type="text/css" href("~/Content/style.css">
13  <link rel="stylesheet" type="text/css" href "~/Content/bootstrap.css">
14  <link rel="stylesheet" type="text/css" href "~/Content/bootstrap.min.css">
15
16  <br/>
17  <hr/>
18  <div class="row" style="display: flex; flex-direction: row; flex-wrap: nowrap; justify-content: space-between;">
19      @foreach (var cats in Model.recommendations)
20      {
21          <div class="gallery_box" style="width: 300px; height: 300px;">
22              
23              <h1>
24                  @Html.ActionLink(cats.ToString(), "Review", "Home", routeValues: new { hotelId = cats, path = p })
25              </h1>
26          </div>
27      }
28  </div>
29
30  <h3>
31      Please click on above hotel links to explore more information about the hotel
32  </h3>
33

```

Solution Explorer

Search Solution Explorer (Ctrl+.)

- ExternalLoginConfirmation.cs
- ExternalLoginFailure.cshtml
- ForgotPassword.cshtml
- ForgotPasswordConfirmation.cs
- Login.cshtml
- Register.cshtml
- ResetPassword.cshtml
- ResetPasswordConfirmation.cs
- SendCode.cshtml
- VerifyCode.cshtml
- Admin
- Home
 - About.cshtml
 - Business.cshtml
 - Cluster.cshtml
 - Contact.cshtml
 - Index.cshtml
 - Report.cshtml
 - Result.cshtml
 - Review.cshtml
 - Visualize.cshtml
- Manage
- Shared

HotelSystem - Microsoft Visual Studio

File Edit View Project Build Debug Team Tools Test Analyze Window Help

Debug Any CPU Google Chrome

Server Explorer Toolbox Cloud Explorer

Cluster.cshtml Result.cshtml Index.cshtml _Layout.cshtml Review.cshtml Classify.cshtml Business.cshtml

```

16      <p>We are the fastest growing hotel recommendation organization. If you are our customer, we value your
17      <p>Let's go big !!</p>
18  </div>
19  <div class="cleaner_h40">&nbsp;</div>
20  <div class="row" style="border: double;">
21      <div class="col-md-6">
22          @using (Html.BeginForm("Business", "Home", FormMethod.Post, new { @class = "form-horizontal", role = "form" }))
23          {
24              @Html.AntiForgeryToken()
25              <h4>Enter values to predict Hotel Category for User</h4>
26              <hr />
27              @Html.ValidationSummary("", new { @class = "text-danger" })
28              <div class="form-group">
29                  @Html.LabelFor(m => m.user_id, new { @class = "col-md-4 control-label" })
30                  <div class="col-md-8">
31                      @Html.DropDownListFor(m => m.user_id, items, new { @class = "form-control" })
32                  </div>
33              </div>
34
35              <div class="form-group">
36                  @Html.LabelFor(m => m.user_location_city, new { @class = "col-md-4 control-label" })
37                  <div class="col-md-8">
38                      @Html.TextBoxFor(m => m.user_location_city, new { @class = "form-control" })
39                  </div>
40
41              <div class="form-group">
42                  @Html.LabelFor(m => m.srch_adults_cnt, new { @class = "col-md-4 control-label" })
43                  <div class="col-md-8">
44

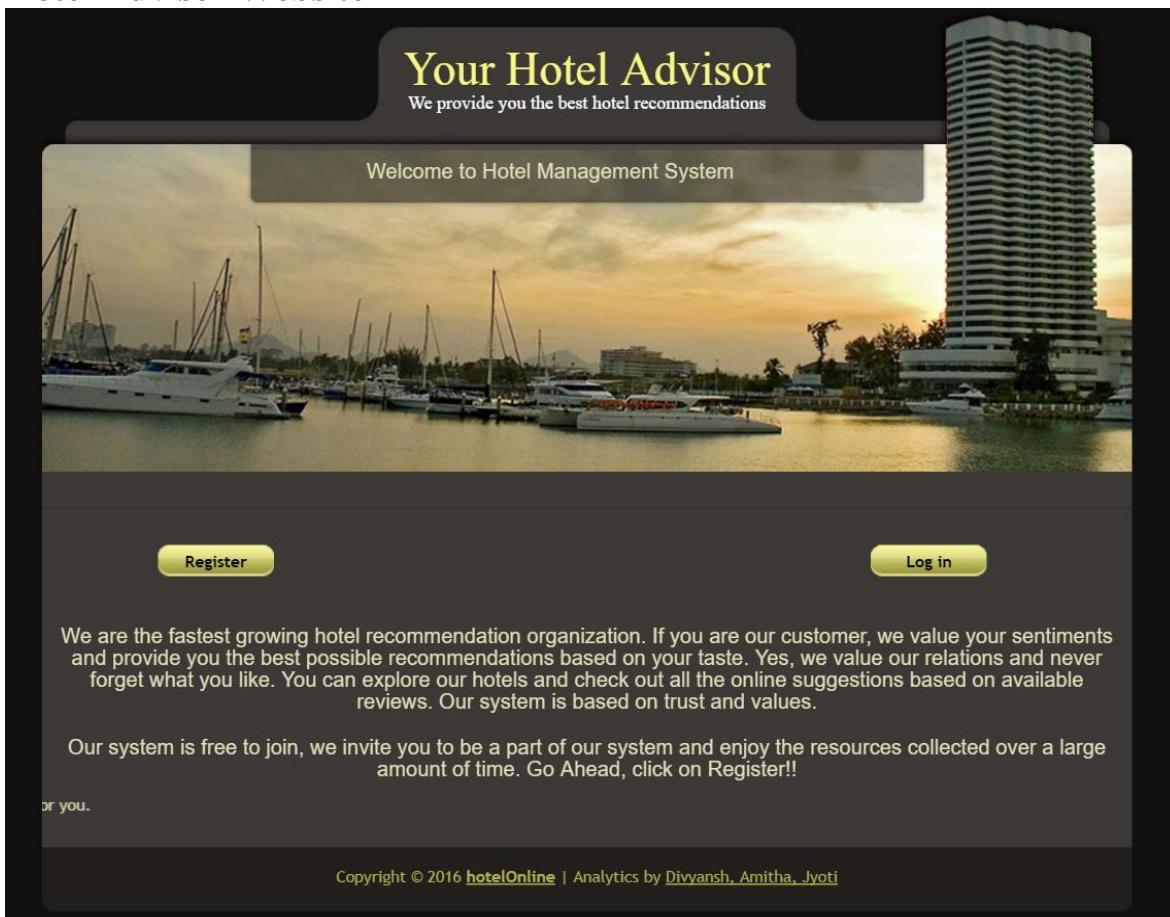
```

Solution Explorer

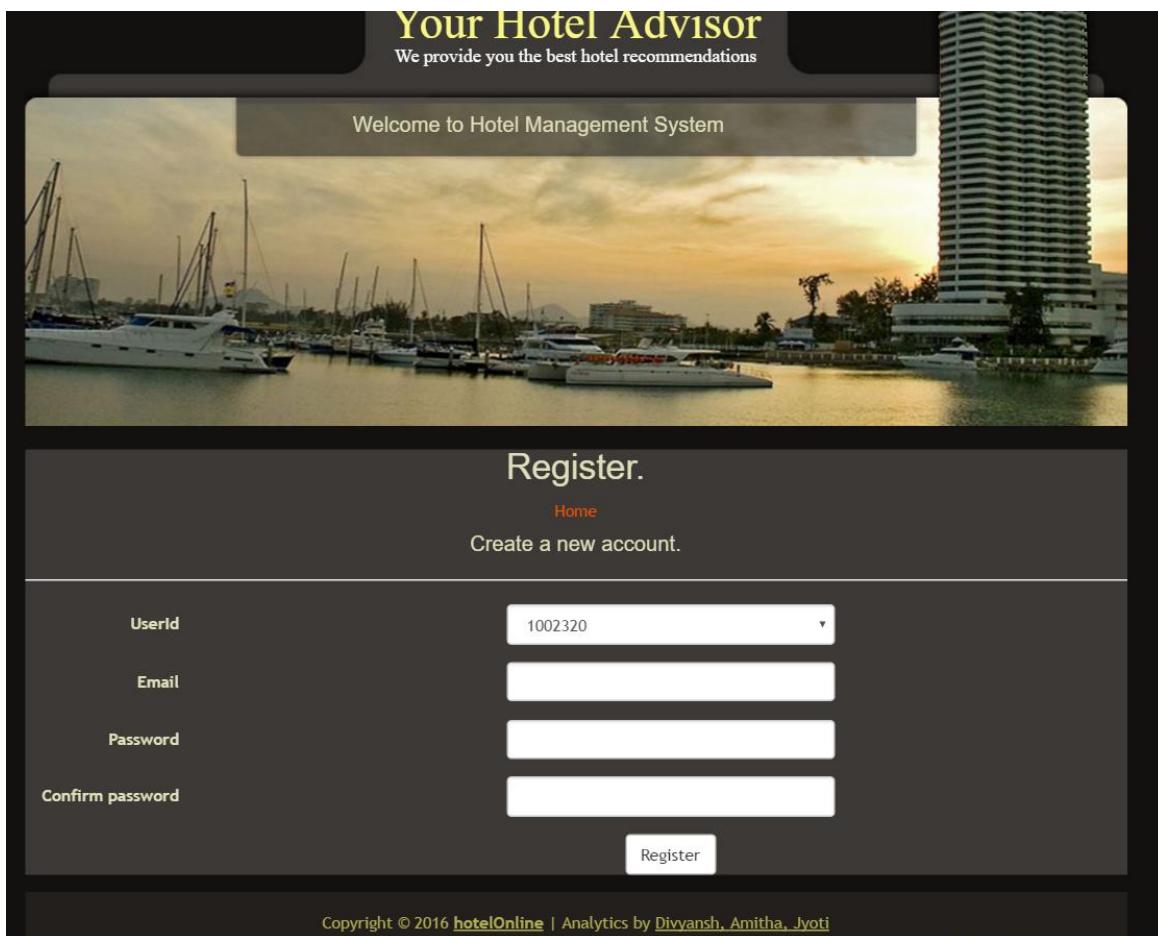
Search Solution Explorer (Ctrl+.)

- ExternalLoginConfirmation.cs
- ExternalLoginFailure.cshtml
- ForgotPassword.cshtml
- ForgotPasswordConfirmation.cs
- Login.cshtml
- Register.cshtml
- ResetPassword.cshtml
- ResetPasswordConfirmation.cs
- SendCode.cshtml
- VerifyCode.cshtml
- Admin
- Home
 - About.cshtml
 - Business.cshtml
 - Cluster.cshtml
 - Contact.cshtml
 - Index.cshtml
 - Report.cshtml
 - Result.cshtml
 - Review.cshtml
 - Visualize.cshtml
- Manage
- Shared

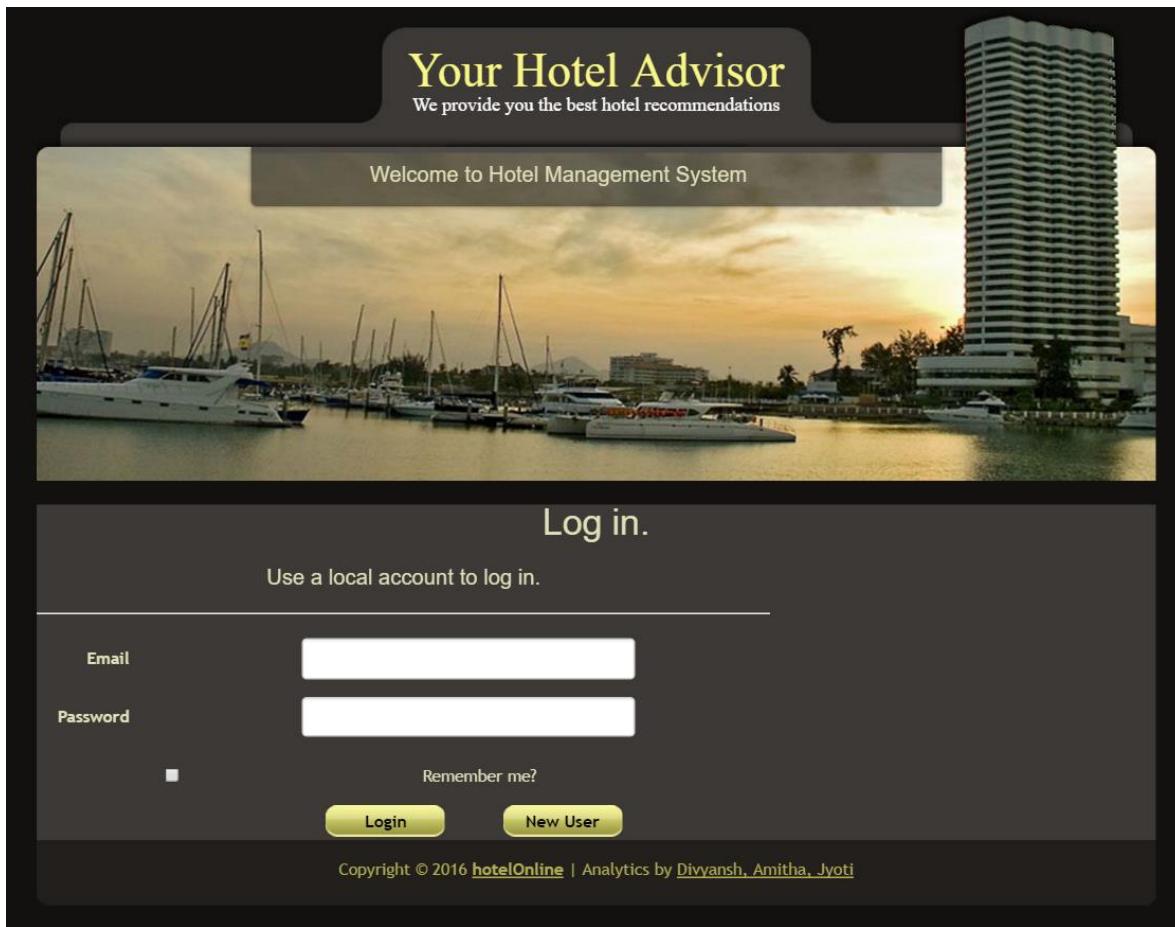
✓ Hotel Advisor Website



✓ Register page



✓ Sign-in page



✓ Admin home page – predict hotel cluster

The screenshot shows an "Admin home page" for predicting hotel clusters. On the left, a form titled "Enter values to predict Hotel Category for User" contains fields for "User ID" (set to 1002320), "User Location City", "Adults(18+)", "Children(0-17)", "Rooms", "Destination Id", "Check in Date" (mm/dd/yyyy), and "Check out Date" (mm/dd/yyyy). Below these fields is a yellow button labeled "Show me cluster". On the right, the heading "Hotel Prediction" is displayed above the text "Cluster Category to be recommended to user :".

User case Data:

Your Hotel Advisor
We provide you the best hotel recommendations

Home Report Visualizations Explore Sample Data Contact

User ID	User Location City	Adults(18+)	Children(0-17)	Rooms	Destination ID	Check in Date	Check out Date
307065	17494	2	0	1	11835	6/27/2017 12:00:00 AM	6/30/2017 12:00:00 AM
194208	48262	2	0	1	8855	1/10/2017 12:00:00 AM	1/17/2017 12:00:00 AM
613844	27731	2	0	1	8744	5/9/2017 12:00:00 AM	5/18/2017 12:00:00 AM
1053478	53517	2	0	1	11353	8/13/2017 12:00:00 AM	8/17/2017 12:00:00 AM
936906	36086	3	0	3	5405	4/4/2017 12:00:00 AM	4/8/2017 12:00:00 AM
653198	47997	1	0	1	8260	10/6/2016 12:00:00 AM	10/10/2016 12:00:00 AM
154070	27102	2	2	2	8740	11/27/2016 12:00:00 AM	11/28/2016 12:00:00 AM
806146	4687	3	1	2	17520	6/18/2017 12:00:00 AM	6/19/2017 12:00:00 AM
720068	34868	5	0	2	23820	4/28/2017 12:00:00 AM	5/10/2017 12:00:00 AM
257191	41779	2	5	1	11939	12/13/2016 12:00:00 AM	12/16/2016 12:00:00 AM
1155959	53078	3	0	2	8289	8/19/2017 12:00:00 AM	8/27/2017 12:00:00 AM

Copyright © 2016 [hotelOnline](#) | Analytics by [Divyansh, Amitha, Jyoti](#)

Recommended hotels for user

Hotel Recommendations for User 1002320



67657 11905 130585

Please click on above hotel links to explore more information about the hotel

hotel review sentiment analysis

Hotel Recommendations for User 1002320

Review

Our analysis on Hotel 11905 :

Visitors Ratings and Reviews: No Review



Hotel ID
Hotel 18735 ▾

Check Hotel Reviews

Visitors Ratings and Reviews:

Excellent

- ✓ Customer home page

1.8 References

1. <http://www.text-analytics101.com/2014/10/computing-precision-and-recall-for.html>
Computing Precision and Recall for Multi-Class Classification Problems
2. <https://github.com/oliviak/Recommender-in-Azure/tree/master/4%20Content-Filtering%20and%20Hybrid%20recommender>
Content Based and Hybrid Recommendation System
3. <https://www.r-bloggers.com/sentiment-analysis-on-donald-trump-using-r-and-tableau/>
Sentiment Analysis in R studio