

October 4, 2020

Hotel booking demand

Inturduction

Here is a data set about hotel booking demand which I would like to explore and to use machine learning methodology. This data set includes features about the type of hotel booked, the lead time before the booking, the number of stays and etc.

This is the order of my work:

- First exploring
- Preprocessing
- Supervised learning
- Unsupervised learning

In [2]:

```
import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
```

Let's read the data set.

In [3]:

```
try:
    hotel_data = pd.read_csv('C:/Users/Matan/Documents/Python/hotel_bookings.csv')
except: print("Something got wrong")
```

In [4]:

```
hotel_data
```

Out [4]:

	hotel	is_canceled	lead_time	arrival_date_year	arrival_date_month	arrival_date_week_number	arrival_date_day_of_month	stays_in_weekend_nights	stays_in_week_nights
0	Resort Hotel	0	342	2015	July	27	1	1	
1	Resort Hotel	0	737	2015	July	27	1	1	
2	Resort Hotel	0	7	2015	July	27	1	1	
3	Resort Hotel	0	13	2015	July	27	1	1	
4	Resort Hotel	0	14	2015	July	27	1	1	
...
119385	City Hotel	0	23	2017	August	35	30		
119386	City Hotel	0	102	2017	August	35	31		
119387	City Hotel	0	34	2017	August	35	31		
119388	City Hotel	0	109	2017	August	35	31		
119389	City Hotel	0	205	2017	August	35	29		

119390 rows x 32 columns

The data is ready.

First exploring

In [5]:

```
print("The number number of rows in hotel data are", hotel_data.shape[0])
print("The number of columns in hotel data are", hotel_data.shape[1])
```

The number number of rows in hotel data are 119390
The number of columns in hotel data are 32

In [6]:

```
hotel_data.columns
```

Out [6]:

```
Index(['hotel', 'is_canceled', 'lead_time', 'arrival_date_year',
       'arrival_date_month', 'arrival_date_week_number',
       'arrival_date_day_of_month', 'stays_in_weekend_nights',
       'stays_in_week_nights', 'adults', 'children', 'babies', 'meal',
       'country', 'market_segment', 'distribution_channel',
       'is_repeated_guest', 'previous_cancellations',
       'previous_bookings_not_canceled', 'reserved_room_type',
       'assigned_room_type', 'booking_changes', 'agent',
       'company', 'days_in_waiting_list', 'customer_type', 'adr',
       'required_car_parking_spaces', 'total_of_special_requests',
       'reservation_status', 'reservation_status_date'],
      dtype='object')
```

In [7]:

```
hotel_data.dtypes
```

Out [7]:

```
hotel                object
is_canceled          int64
lead_time            int64
arrival_date_year    int64
arrival_date_month   int64
arrival_date_week_number int64
arrival_date_day_of_month int64
stays_in_weekend_nights int64
stays_in_week_nights int64
adults              int64
children            float64
babies              int64
meal                object
country             object
market_segment      object
distribution_channel object
is_repeated_guest   int64
previous_cancellations object
previous_bookings_not_canceled int64
reserved_room_type  object
assigned_room_type  object
booking_changes     int64
agent              float64
company            float64
days_in_waiting_list int64
customer_type       object
required_car_parking_spaces int64
total_of_special_requests int64
reservation_status   object
reservation_status_date object
dtype: object
```

In [8]:

```
hotel_data['country'].drop_duplicates().count()
```

Out [8]:

```
177
```

In [9]:

```
hotel_data['distribution_channel'].drop_duplicates().count()
```

Out [9]:

```
5
```

In [10]:

```
hotel_data['hotel'].drop_duplicates()
```

Out [10]:

```
0    Resort Hotel
40060    City Hotel
Name: hotel, dtype: object
```

There are only two types of hotels.

In [11]:

```
hotel_data['hotel'].value_counts().plot(kind='bar', title="Number of Resort and City Hotels")
plt.grid()
```



The number of city hotels is twice as large as in resort hotel in this data set.

In [12]:

```
hotel_kinds = hotel_data.groupby('hotel')
for name, kind in hotel_kinds:
    print("The average number of children and babies in", name, "are:", kind['children'].mean(), kind['babies'].mean())
```

The average number of children and babies in City Hotel are: 0.09136799048483473 0.004941384091768561
The average number of children and babies in Resort Hotel are: 0.128681977033446834 0.0139041437843235

We can learn from this result that children and babies are more likely to go to resort hotels.

Now let's check the distribution of the number of days that elapsed between the entering date of the booking and the arrival date.

In [13]:

```
sns.distplot(hotel_data['lead_time'], bins=20)
plt.grid()
plt.xlabel('Lead Time')
plt.ylabel('Count')
plt.title("The Distribution Of The Days Before The Vacation")
```

Out [13]:

```
Text(0.5, 1.0, 'The Distribution Of The Days Before The Vacation')
```



We can learn from this graph that most of people choose to book their hotel up to 100 days from their arrival date.

Let's check which country consists the most hotels in this data set.

In [14]:

```
countries = hotel_data.groupby('country')['country'].value_counts().head(10)
countries.plot(kind='bar', color='m', linewidth=2)
plt.xlabel('Countries')
plt.ylabel('Count')
plt.title("The number of hotels in each country")
```

Out [14]:

```
Text(0.5, 1.0, 'The number of hotels in each country')
```



We can see that "AGG" and "ARG" consists of a lot more hotels than the others in this data set. Moreover, the average number of hotels in each country in this data set is a bit less than 50.

Checking the numbers of repeated guests in each type of hotel.

In [15]:

```
hotel_data.groupby('hotel')['is_repeated_guest'].sum().plot(kind='bar', color='royalblue', alpha=0.3)
plt.grid(color='black', linestyle='--', linewidth=1, axis='y', alpha=0.9)
plt.title("The number of repeated guests in each type of hotel")
```

Out [15]:

```
Text(0.5, 1.0, 'The number of repeated guests in each type of hotel')
```



The city hotels consists of a bit more repeated guests than the resort hotel. However, the gap between those kinds of hotel is smaller than the gap of the number of the hotels in each type. we can infer that guests are more likely visit again in resort hotels and not in city hotels.

Checking the number of adults in each type of hotel.

In [16]:

```
hotel_data.groupby('hotel')['adults'].sum().plot(kind='bar', color='royalblue', alpha=0.3)
plt.grid(color='black', linestyle='--', linewidth=1, axis='y', alpha=0.9)
graph = sns.FacetGrid(hotel_data, col='hotel')
graph.map(plt.hist, 'lead_time', bins=20)
```

Out [16]:

```
Text(0.5, 1.0, 'The number of adults in each type of hotel')
```



The relative gap is pretty the same as the relative gap in the number of hotels consists of the data set.

Now let's see if there is a difference between the lead time in city hotels to resort hotels.

In [17]:

```
graph = sns.FacetGrid(hotel_data, col='hotel')
seaborn.axisgrid.FacetGrid at 0x20ed56480
```

Out [17]:



Both are pretty the same.

Checking NaN values.

In [18]:

```
hotel_data.isnull().sum()
```

Out [18]:

```
hotel                0
is_canceled          0
lead_time            0
arrival_date_year    0
arrival_date_month   0
arrival_date_week_number 0
arrival_date_day_of_month 0
stays_in_weekend_nights 0
stays_in_week_nights 0
adults              4
children            0
babies              0
meal                498
country             0
market_segment      0
distribution_channel 0
is_repeated_guest   0
previous_cancellations 0
previous_bookings_not_canceled 0
reserved_room_type  0
assigned_room_type  0
booking_changes     0
agent              0
company            16340
adr                112593
days_in_waiting_list 0
customer_type       0
required_car_parking_spaces 0
total_of_special_requests 0
reservation_status   0
reservation_status_date 0
dtype: int64
```

Most of the values in the company column are missing.
13.68% of the values in the agent column are missing as well.

Now let's check if there is correlation between the features.

In [81]:

```
sns.heatmap(hotel_data.corr(), annot = True)
plt.title('Heatmap', fontsize = 20)
plt.rcParams['figure.dpi'] = [16, 12]
graph = sns.FacetGrid(hotel_data, col='hotel')
plt.rcParams['figure.dpi'] = 200
```



The correlation of most of the features is around 0. However, there is a positive correlation (0.5) between the stays in a week to weekend nights and the market segment to distribution (0.77).
In addition, there is a negative correlation (-0.53) between the arrival date year to the month.

After I explored the data, my two main questions are:

- The question for the supervised learning is: Can we predict a type of hotel by the feature of booking?

- The question for the unsupervised learning is: Can we split our data to N clusters such that every cluster is including similar booking features?

Preprocessing

In [19]:

```
from sklearn import preprocessing
```

First, I will delete unnecessary columns and rows which consists of NaN values.

In [20]:

```
hotel_data = hotel_data.drop(columns=['arrival_date_week_number', 'country', 'reserved_room_type', 'assigned_room_type', 'agent', 'company', 'reservation_status_date'])
hotel_data = hotel_data.dropna(subset=['children'])
```

Deleting explanation

- Columns:
- Arrival date week number - Two many different values that couldn't help to solve our main questions.
 - Country - Two many different countries which can cause us to overfitting.
 - Reserved and assigned room type - not belong to our main questions.
 - Agent - Missing values.
 - Company - Missing values.
 - Reservation status date - not belong to our main questions.
- Changing the name of the months to their numbers.
- In [21]:
- ```
months = {'Jan': 1, 'Feb': 2, 'Mar': 3, 'Apr': 4, 'May': 5, 'Jun': 6, 'Jul': 7, 'Aug': 8, 'Sep': 9, 'Oct': 10, 'Nov': 11, 'Dec': 12}
hotel_data.arrival_date_month = hotel_data.arrival_date_month.map(months)
```

Now I will do the preprocessing. One data will fit supervised learning and the other to unsupervised.

In [23]:

```
def first_preprocessing(data_frame, supervised_or_unsupervised):
 try:
 new_data = data_frame
 if supervised_or_unsupervised == 'supervised':
 encoders = {'hotel': preprocessing.LabelEncoder(),
 'meal': preprocessing.LabelEncoder(),
 'market_segment': preprocessing.LabelEncoder(),
 'distribution_channel': preprocessing.LabelEncoder(),
 'deposit_type': preprocessing.LabelEncoder(),
 'customer_type': preprocessing.LabelEncoder(),
 'reservation_status': preprocessing.LabelEncoder()}
 encoder_columns_list = ['hotel', 'meal', 'market_segment', 'distribution_channel', 'deposit_type', 'customer_type', 'reservation_status']
 for column in encoder_columns_list:
 new_data[column] = encoders[column].fit_transform(new_data[column].astype(str))
 elif supervised_or_unsupervised == 'unsupervised':
 new_data = pd.get_dummies(new_data)
 columns = new_data.columns
 min_max_scaler = preprocessing.MinMaxScaler()
 new_data = min_max_scaler.fit_transform(new_data)
 new_data = pd.DataFrame(new_data)
 new_data.columns = columns
 except: print("Something got wrong - first_preprocessing")
```

Preprocessing explanation:

Supervised data:

- Changing the columns which contain object to int values. #### Unsupervised data:
- Using "One Hot Vector" to the columns that contain object values.
- Normalize the whole data set.

In [24]:

```
hotel_supervised = first_preprocessing(hotel_data, 'supervised')
hotel_supervised
```

Out [24]:

|        | hotel | is_canceled | lead_time | arrival_date_year | arrival_date_month | arrival_date_day_of_month | stays_in_weekend_nights | stays_in_week_nights |
|--------|-------|-------------|-----------|-------------------|--------------------|---------------------------|-------------------------|----------------------|
| 0      | 1     | 0           | 342       | 2015              | 7                  | 1                         | 1                       | 0                    |
| 1      | 1     | 0           | 737       | 2015              | 7                  | 1                         | 1                       | 0                    |
| 2      | 1     | 0           | 7         | 2015              | 7                  | 1                         | 1                       | 0                    |
| 3      | 1     | 0           | 13        | 2015              | 7                  | 1                         | 1                       | 0                    |
| 4      | 1     | 0           | 14        | 2015              | 7                  | 1                         | 1                       | 0                    |
| ...    | ...   | ...         | ...       | ...               | ...                | ...                       | ...                     | ...                  |
| 119385 | 0     | 0           | 23        | 2017              | 8                  | 30                        | 3                       | 2                    |
| 119386 | 0     | 0           | 102       | 2017              | 8                  | 31                        | 2                       | 2                    |
| 119387 | 0     | 0           | 34        | 2017              | 8                  | 31                        | 2                       | 2                    |
| 119388 | 0     | 0           | 109       | 2017              | 8                  | 31                        | 2                       | 2                    |
| 119389 | 0     | 0           | 205       | 2017              | 8                  | 29                        | 2                       | 2                    |

119390 rows x 25 columns

Supervised data is ready.

In [25]:

```
hotel_unsupervised = first_preprocessing(hotel_data, 'unsupervised')
```

Out [25]:

|        | hotel | is_canceled | lead_time | arrival_date_year | arrival_date_month | arrival_date_day_of_month | stays_in_weekend_nights | stays_in_week_nights |
|--------|-------|-------------|-----------|-------------------|--------------------|---------------------------|-------------------------|----------------------|
| 0      | 1     | 0           | 0.464043  | 0.0               | 0.545455           | 0.000000                  | 0.000000                | 0.000000             |
| 1      | 1     | 0           | 1.000000  | 0.0               | 0.545455           | 0.000000                  | 0.000000                | 0.000000             |
| 2      | 1     | 0           | 0.009498  | 0.0               | 0.545455           | 0.000000                  | 0.000000                | 0.000000             |
| 3      | 1     | 0           | 0.017639  | 0.0               | 0.545455           | 0.000000                  | 0.000000                | 0.000000             |
| 4      | 1     | 0           | 0.018896  | 0.0               | 0.545455           | 0.000000                  | 0.000000                | 0.000000             |
| ...    | ...   | ...         | ...       | ...               | ...                | ...                       | ...                     | ...                  |
| 119381 | 0.0   | 0.0         | 0.031208  | 1.0               | 0.636364           | 0.966667                  | 0.105263                | 0.105263             |
| 119382 | 0.0   | 0.0         | 0.138399  | 1.0               | 0.636364           | 1.000000                  | 0.105263                | 0.105263             |
| 119383 | 0.0   | 0.0         | 0.046133  | 1.0               | 0.636364           | 1.000000                  | 0.105263                | 0.105263             |
| 119384 | 0.0   | 0.0         | 0.174897  | 1.0               | 0.636364           | 1.000000                  | 0.105263                | 0.105263             |
| 119385 | 0.0   | 0.0         | 0.278155  | 1.0               | 0.636364           | 0.933333                  | 0.105263                | 0.105263             |

119390 rows x 25 columns

Unsupervised data is ready.

Now we can predict a type of hotel by the feature of booking?

I will use supervised learning with the help of Naive Bayes Classifier. This classifier is a statistical classification technique based on Bayes Theorem.

Following the main question, I will try to predict the right type of hotel between the two:

- Resort hotel
- City hotel

In [32]:

```
from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import confusion_matrix, confusion_matrix, precision_score, recall_score
```

Separating the target column and all the rest features:

In [33]:

```
X = hotel_supervised.drop(columns = ['hotel']).copy()
y = hotel_supervised['hotel']
```

Building two functions that supposed to split the data for train and test and create the Naive Bayes Classifier:

In [34]:

```
def split_test_train(X, y, test_size):
 try:
 return train_test_split(X, y, test_size=test_size, random_state=0)
 except:
 print("Something got wrong - split_test_train")
```

In [35]:

```
def create_naive_bayes_classifier(X, y):
 try:
 model = GaussianNB()
 model.fit(X, y)
 return model
 except:
 print("Something got wrong - create_naive_bayes_classifier")
```

Before I choose the train and the test part, I want to make a plot that can show me what is the best accuracy based on the partition of the train and test.

In [36]:

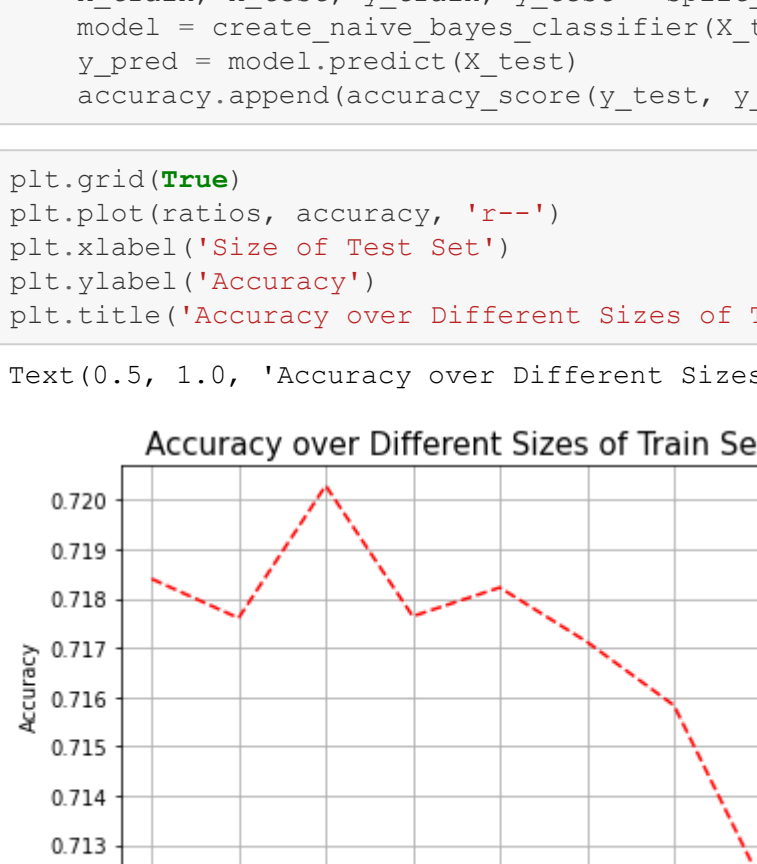
```
ratios = np.arange(0.1, 0.9, 0.1)
accuracy = []
for ratio in ratios:
 X_train, X_test, y_train, y_test = split_test_train(X, y, ratio)
 model = create_naive_bayes_classifier(X_train, y_train)
 y_pred = model.predict(X_test)
 accuracy.append(accuracy_score(y_test, y_pred))
```

In [37]:

```
plt.plot(ratios, accuracy, 'r--')
plt.xlabel('Size of Test Set')
plt.ylabel('Accuracy')
plt.title('Accuracy Over Different Sizes of Train Set', fontsize=15)
```

Out [37]:

```
Text(0.5, 1.0, 'Accuracy Over Different Sizes of Train Set')
```



As we can see, The best accuracy I will get is if I will separate the data to 0.3 test group and 0.7 train group.

Now we are ready to activate the algorithm and analyze the results.

In [38]:

```
X_train, X_test, y_train, y_test = split_test_train(X, y, 0.3)
```

In [39]:

```
model = create_naive_bayes_classifier(X_train, y_train)
```

In [40]:

```
y_pred = model.predict(X_test)
```

In [41]:

```
accuracy_score(y_test, y_pred)
```

Out [41]:

```
0.702920606562431
```

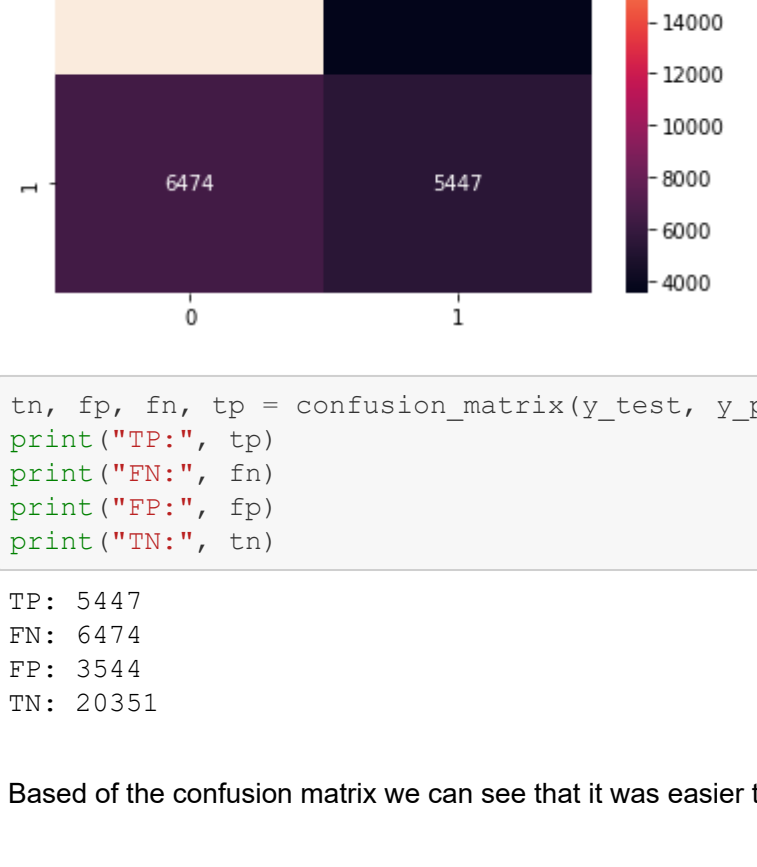
We received not bad accuracy depends on the starting point which our chances were 50/50.

Let's use a confusion matrix to see if we can find something significant.

In [42]:

```
sns.heatmap(confusion_matrix(y_test, y_pred), annot=True, fmt="d")
```

Out [42]:



In [43]:

```
tn, fp, fn, tp = confusion_matrix(y_test, y_pred).ravel()
print("TN:", tn)
print("FP:", fp)
print("FN:", fn)
print("TP:", tp)
```

TN: 20351  
FP: 6474  
FN: 3544  
TP: 3447

Based on the confusion matrix we can see that it was easier to the algorithm to predict the city hotel(0).

To dig more into the result, I will use two rates:

- Precision - this rate using all the points predicted to positive to how many of them are actually positive. In our cases, the positive is the resort hotel.
- Recall - this is the ratio of the number correctly predicted positive class (resort hotel) to the total number of positive classes.

In [44]:

```
precision_score(y_test, y_pred)
```

Out [44]:

```
0.6058280502724948
```

From the precision score, we learn that when the algorithm predicts resort hotel, in most of the cases he was right (60.5%).

In [45]:

```
recall_score(y_test, y_pred)
```

Out [45]:

```
0.456924756346783
```

From the recall score we learn that most of the predictions of the actual resort hotel was false (45.6%).

Below are two plots about the difference between the actual data to the predicted one.

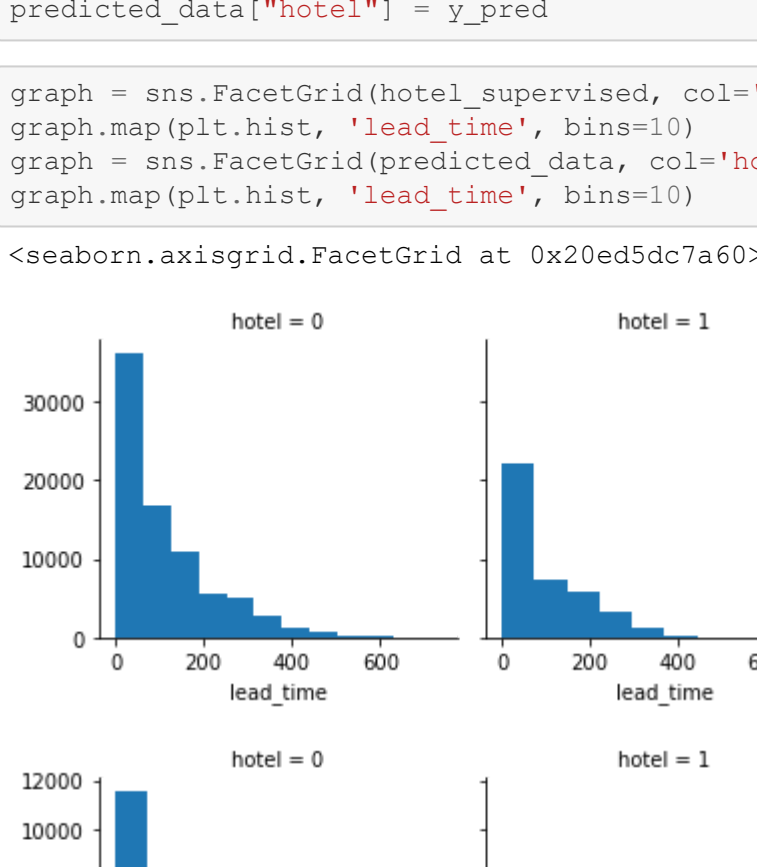
In [46]:

```
predicted_data = X_test.copy()
predicted_data["hotel"] = y_pred
```

In [47]:

```
graph = sns.FacetGrid(hotel_supervised, col='hotel')
graph.map(plt.hist, 'lead_time', bins=10)
graph = sns.FacetGrid(predicted_data, col='hotel')
graph.map(plt.hist, 'lead_time', bins=10)
```

Out [47]:

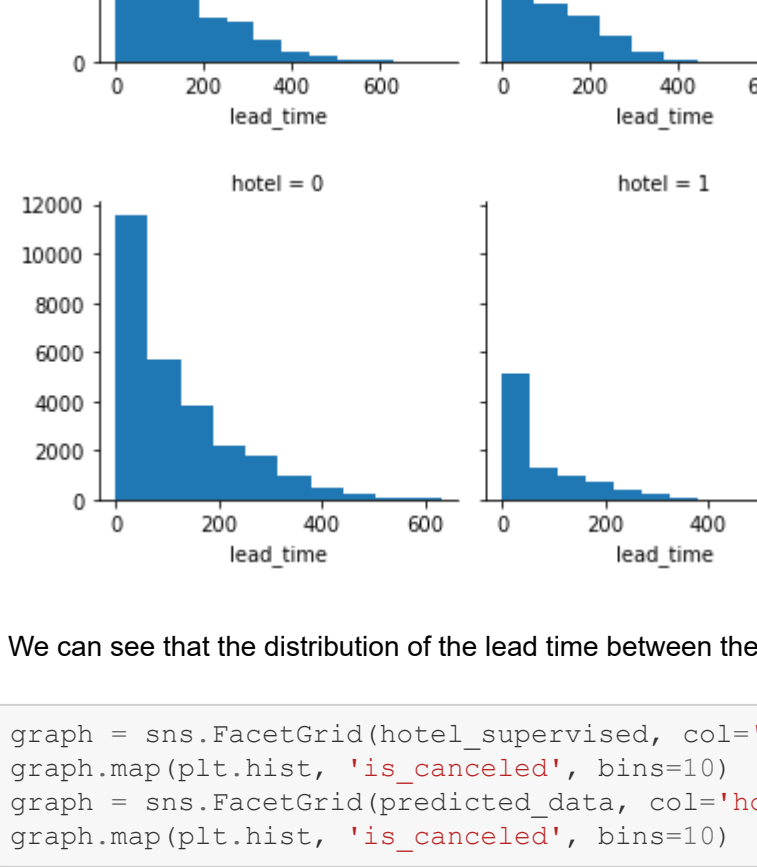


We can see that the distribution of the lead time between the actual data to the predicted data is pretty the same.

In [48]:

```
graph = sns.FacetGrid(hotel_supervised, col='hotel')
graph.map(plt.hist, 'is_canceled', bins=10)
graph = sns.FacetGrid(predicted_data, col='hotel')
graph.map(plt.hist, 'is_canceled', bins=10)
```

Out [48]:



The relative gap between the cancellation of the actual data is a bit less than the predicted one.

After I got the result and I analyzed it, I wanted to check if I can optimize the algorithm by keeping the features that clarify the most the prediction (city/resort hotel). Therefore, I am going to use two methods about feature selection and after I will keep only the most important features. I will activate the algorithm one more time to check if the prediction was more accurate.

In [49]:

```
from sklearn.feature_selection import RandomForestClassifier
I will use supervised learning with the help of Naive Bayes Classifier. This classifier is a statistical classification technique based on Bayes Theorem.
```

Method 1 - Tree based Feature Selection and Random Forest Classification

In this method, I will use Random Forest algorithm to predict the type of hotel. After that, I will use 'feature/importances' to check the contribution of any feature to the prediction and in the end, I will display it on bar plot.

In [50]:

```
clf_rf = RandomForestClassifier(n_estimators=20)
clf_rf = clf_rf.fit(X_train, y_train)
importances = clf_rf.feature_importances_
std = np.std([tree.feature_importances_ for tree in clf_rf.estimators_], axis=0)
indices = np.argsort(importances)[::-1]
```



```
[51]: print("Feature ranking:")
for f in range(X_train.shape[1]):
 print("%d. feature %d (%2f)" % (f + 1, indices[f], importances[indices[f]]))
```

```
Feature ranking:
1. feature 20 (0.303922)
2. feature 8 (0.190933)
3. feature 1 (0.089587)
4. feature 10 (0.074752)
5. feature 3 (0.072111)
6. feature 4 (0.060585)
7. feature 21 (0.038973)
8. feature 11 (0.037115)
9. feature 5 (0.031120)
10. feature 2 (0.027636)
11. feature 7 (0.023431)
12. feature 19 (0.021380)
13. feature 22 (0.021327)
14. feature 12 (0.021293)
15. feature 8 (0.013677)
16. feature 16 (0.012840)
17. feature 17 (0.009063)
18. feature 0 (0.008365)
19. feature 23 (0.005573)
20. feature 15 (0.004977)
21. feature 14 (0.004363)
22. feature 13 (0.002232)
23. feature 18 (0.002927)
24. feature 9 (0.001552)
```



We can learn that the feature "ad" is significantly more clearly to the prediction in comparison to the other features, for instance, this feature almost 3 times bigger (2.877) than the second most important feature.

## Method 2 - Recursive feature elimination (RFE) with Random Forest Classification

The algorithm assigns weights to each of features. Whose absolute weights are the smallest are pruned from the current set features. That procedure is recursively repeated on the pruned set until the desired number of features.

In [53]:

```
clf_rf = RandomForestClassifier(n_estimators=20)
rfe = RFE(estimator=clf_rf, n_features_to_select=11, step=1)
rfe = rfe.fit(X_train, y_train)
```

In [54]:

```
print("Chosen best 11 feature by RFE:", X_train.columns[rfe.support_])
```

Chosen best 11 feature by RFE: Index(['lead\_time', 'arrival\_date\_year', 'arrival\_date\_month', 'arrival\_date\_day\_of\_month', 'stays\_in\_weekend\_nights', 'stays\_in\_week\_nights', 'meal', 'market\_segment', 'distribution\_channel', 'adr', 'required\_car\_parking\_spaces'], dtype='object')

This method brought me the same 11 most important features as the first method, therefore I will use only those 11 features and will activate the algorithm once again.

In [55]:

```
new_hotel_supervised = hotel_supervised.drop(columns = ['deposit_type', 'children', 'booking_changes', 'reservation_status', 'is_canceled', 'previous_bookings_not_canceled', 'previous_cancellations', 'days_in_waiting_list', 'is_repeated_guest', 'babies', 'distribution_channel', 'deposit_type'])
```

In [56]:

```
X = new_hotel_supervised.drop(columns = ['hotel']).copy()
y = new_hotel_supervised['hotel'].copy()
```

In [57]:

```
X_train, X_test, y_train, y_test = split_test_train(X, y, 0.3)
```

In [58]:

```
model = create_naive_bayes_classifier(X_train, y_train)
```

In [59]:

```
y_pred = model.predict(X_test)
```

In [60]:

```
accuracy_score(y_test, y_pred)
```

Out[60]: 0.7284921152557516

I increased the accuracy by 0.008 points - unfortunately not so significantly.

In [61]:

```
sns.heatmap(confusion_matrix(y_test, y_pred), annot=True, fmt="d")
```

Out[61]:

From the confusion matrix, we can see that the algorithm predicted better the city hotel. However, the prediction of the resulting hotel was less good than before.

In [63]:

```
hotel_type = hotel_supervised.groupby('hotel')[['lead_time', 'arrival_date_year', 'arrival_date_month', 'arrival_date_day_of_month', 'stays_in_weekend_nights', 'stays_in_week_nights', 'meal', 'market_segment', 'distribution_channel', 'adr', 'required_car_parking_spaces', 'total_of_special_requests']].mean()
```

Out[63]:

| hotel | lead_time  | arrival_date_year | arrival_date_month | arrival_date_day_of_month | stays_in_weekend_nights | stays_in_week_nights | adr    |
|-------|------------|-------------------|--------------------|---------------------------|-------------------------|----------------------|--------|
| 0     | 109.741108 | 2016.174344       | 0.856400           | 15.787004                 | 0.706187                | 2.182954             | 1.8501 |
| 1     | 92.675686  | 2016.121443       | 6.544583           | 0.721243                  | 1.189185                | 3.129732             | 1.867  |

In conclusion, in both tests, the accuracy of the algorithm was around 0.72. The feature that has the most clarity about the city of the hotel is the "lead\_time" which is the only rate. The "adr" of the city hotel is significantly higher than the resort hotel (on average of 10.348 higher).

Except this feature, I can say that people book for a longer vacation in resort hotel since the average of the stays in the week and the weekend night is longer in resort hotel than the city hotel. Moreover, the lead time of the booking for the city hotel is longer than the resort hotel (on average of 17.07 longer).

Finally, although the accuracy of the algorithm was above 0.5 (our starting point) the algorithm got difficulties to predict the resort hotel since the number of the city hotel is much higher in this data set than the resort hotel, and the range of almost every feature was very alike to the other type of hotel.

Now let's move one to the next the unsupervised question:

Can we split our data to k clusters such that every cluster is including similar booking features?

To answer that, I will use KMeans algorithm. This algorithm is a method of vector quantization, originally from signal processing, that aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean, serving as a prototype of the cluster.

In [66]:

```
from sklearn.cluster import KMeans
from sklearn.metrics import silhouette_score
```

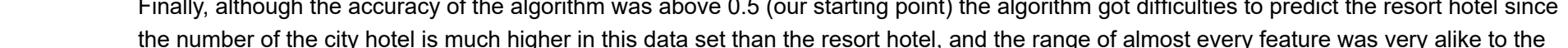
Firstly, I need to find out how many clusters I need to divide into each data I used. In order to deal with this dilemma, I will build a function which called "number\_of\_clusters". This function will return us two graphs. The first graph is called "Sum of square error", in this graph we will see that is the sum of square error as long as the number of clusters increases. Now comes the question of what is the number of clusters that are "suspicious" to be optimal. For answering this question, I will use the "Elbow Method". In this method, when we see a break in some number clusters, we will take that number as an option about how many clusters should we use in this data set. The second graph I will use is called "Silhouette", the silhouette is a measure about how good our algorithm work, his range is between -1 to 1. As long as the value is close to 1 we can say that the clustering was good. In our graph, I will wish to choose the number of clusters by the best silhouette I got. I need to find the best cluster by a combination of those two graphs, a number of clusters which has a good "Elbow break" and ones that has the best silhouette.

In [67]:

```
def number_of_clusters(data_frame):
 try:
 sum_squared = []
 silhouette = []
 for i in range(2, 11):
 kmeans = KMeans(n_clusters=i, init='k-means++')
 kmeans.fit(data_frame)
 sum_squared.append(kmeans.inertia_)
 silhouette.append(silhouette_score(data_frame, kmeans.labels_))
 x1 = (range(2, 11))
 x2 = (range(2, 11))
 y1 = sum_squared
 y2 = silhouette
 plt.subplot(2, 1, 1)
 plt.plot(x1, y1)
 plt.title('Sum of Squared Error (R_2)^2', fontsize=15)
 plt.grid()
 plt.xlabel('R_2^2')
 plt.subplot(2, 1, 2)
 plt.plot(x2, y2)
 plt.title('Silhouette', fontsize=15)
 plt.xlabel('No. of Clusters')
 plt.ylabel('Silhouette')
 plt.grid()
 plt.show()
 except:
 print("Something got wrong - number of clusters")
```

In [68]:

```
number_of_clusters(hotel_unsupervised)
```



As we can see, in the first plot, the sharper elbow break is when I divide the data into 3 clusters. However, when we look at the second plot, the highest silhouette is when I divide the data into 4 clusters. In the second plot, we see that although the best silhouette is to divide the data into 4 clusters, the silhouette of 3 clusters is not much lower than the 4. Moreover, the silhouette score is low in those both option, so I will choose to divide the data into 3 clusters.

Below are two functions:

- create\_kmeans\_classifier - creating KMeans algorithm.
- clusters\_information - display each cluster and his statistic details.

In [69]:

```
def create_kmeans_classifier(k):
 try:
 return KMeans(n_clusters=k, init='k-means++')
 except:
 print("Something got wrong - create_kmeans_classifier")
```

In [70]:

```
def clusters_information(data_frame):
 try:
 clusters = data_frame.groupby("label")
 for name, group in clusters:
 print(name)
 print(group)
 print(group.describe())
 except:
 print("Something got wrong - clusters_information")
```

Activating the functions:

In [71]:

```
kmeans = create_kmeans_classifier(3)
kmeans.fit(hotel_unsupervised)
silhouette_score(hotel_unsupervised, kmeans.labels_)
```

Out[71]: 0.26888207260816356

In [72]:

```
hotel_unsupervised["label"] = pd.Series(kmeans.labels_)
clusters_information(hotel_unsupervised)
```

0

| hotel | is_canceled | lead_time | arrival_date_year | arrival_date_month | \        |
|-------|-------------|-----------|-------------------|--------------------|----------|
| 0     | 1.0         | 0.0       | 0.460403          | 0.0                | 0.545455 |
| 1     | 1.0         | 0.0       | 1.000000          | 0.0                | 0.545455 |
| 2     | 1.0         | 0.0       | 0.009498          | 0.0                | 0.545455 |
| 3     | 1.0         | 0.0       | 0.017639          | 0.0                | 0.545455 |
| 4     | 1.0         | 0.0       | 0.018996          | 0.0                | 0.545455 |

40055 1.0 0.0 0.287653 1.0 0.636364

| hotel | is_canceled | lead_time | arrival_date_year | arrival_date_month | \        |
|-------|-------------|-----------|-------------------|--------------------|----------|
| 40056 | 1.0         | 0.0       | 0.229308          | 1.0                | 0.636364 |
| 40057 | 0.933333    | 0.0       | 0.000000          | 0.210526           | 0.636364 |
| 40058 | 1.000000    | 0.0       | 0.210526          | 0.210526           | 0.636364 |
| 40059 | 1.000000    | 0.0       | 0.210526          | 0.210526           | 0.636364 |

40055 1.0 0.0 0.287653 1.0 0.636364

| hotel | is_canceled | lead_time | arrival_date_year | arrival_date_month | \        |
|-------|-------------|-----------|-------------------|--------------------|----------|
| 40056 | 1.0         | 0.0       | 0.229308          | 1.0                | 0.636364 |
| 40057 | 0.933333    | 0.0       | 0.000000          | 0.210526           | 0.636364 |
| 40058 | 1.000000    | 0.0       | 0.210526          | 0.210526           | 0.636364 |
| 40059 | 1.000000    | 0.0       | 0.210526          | 0.210526           | 0.636364 |

40056 1.0 0.0 0.229308 1.0 0.636364

| hotel | is_canceled | lead_time | arrival_date_year | arrival_date_month | \        |
|-------|-------------|-----------|-------------------|--------------------|----------|
| 40057 | 0.933333    | 0.0       | 0.000000          | 0.210526           | 0.636364 |
| 40058 | 1.000000    | 0.0       | 0.210526          | 0.210526           | 0.636364 |
| 40059 | 1.000000    | 0.0       | 0.210526          | 0.210526           | 0.636364 |

40058 1.0 0.0 0.210526 0.210526 0.636364

| hotel | is_canceled | lead_time | arrival_date_year | arrival_date_month | \        |
|-------|-------------|-----------|-------------------|--------------------|----------|
| 40059 | 1.000000    | 0.0       | 0.210526          | 0.210526           | 0.636364 |

40059 1.0 0.0 0.210526 0.210526 0.636364

| hotel | is_canceled | lead_time | arrival_date_year | arrival_date_month | \        |
|-------|-------------|-----------|-------------------|--------------------|----------|
| 40060 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40061 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40062 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40063 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40064 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40065 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40066 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40067 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40068 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40069 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40070 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40071 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40072 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |

40060 0.0 0.0 0.000000 0.000000 0.000000

| hotel | is_canceled | lead_time | arrival_date_year | arrival_date_month | \        |
|-------|-------------|-----------|-------------------|--------------------|----------|
| 40061 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40062 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40063 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40064 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40065 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40066 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40067 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40068 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40069 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40070 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40071 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40072 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |

40061 0.0 0.0 0.000000 0.000000 0.000000

| hotel | is_canceled | lead_time | arrival_date_year | arrival_date_month | \        |
|-------|-------------|-----------|-------------------|--------------------|----------|
| 40062 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40063 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40064 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40065 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40066 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40067 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40068 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40069 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40070 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40071 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40072 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |

40062 0.0 0.0 0.000000 0.000000 0.000000

| hotel | is_canceled | lead_time | arrival_date_year | arrival_date_month | \        |
|-------|-------------|-----------|-------------------|--------------------|----------|
| 40063 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40064 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40065 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40066 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40067 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40068 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40069 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40070 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40071 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40072 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |

40063 0.0 0.0 0.000000 0.000000 0.000000

| hotel | is_canceled | lead_time | arrival_date_year | arrival_date_month | \        |
|-------|-------------|-----------|-------------------|--------------------|----------|
| 40064 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40065 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40066 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40067 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40068 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40069 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40070 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40071 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40072 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |

40064 0.0 0.0 0.000000 0.000000 0.000000

| hotel | is_canceled | lead_time | arrival_date_year | arrival_date_month | \        |
|-------|-------------|-----------|-------------------|--------------------|----------|
| 40065 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40066 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40067 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40068 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40069 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40070 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40071 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40072 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |

40065 0.0 0.0 0.000000 0.000000 0.000000

| hotel | is_canceled | lead_time | arrival_date_year | arrival_date_month | \        |
|-------|-------------|-----------|-------------------|--------------------|----------|
| 40066 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40067 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40068 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40069 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40070 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40071 | 0.0         | 0.0       | 0.000000          | 0.000000           | 0.000000 |
| 40072 | 0.0         | 0.0       | 0.000000          | 0.000              |          |