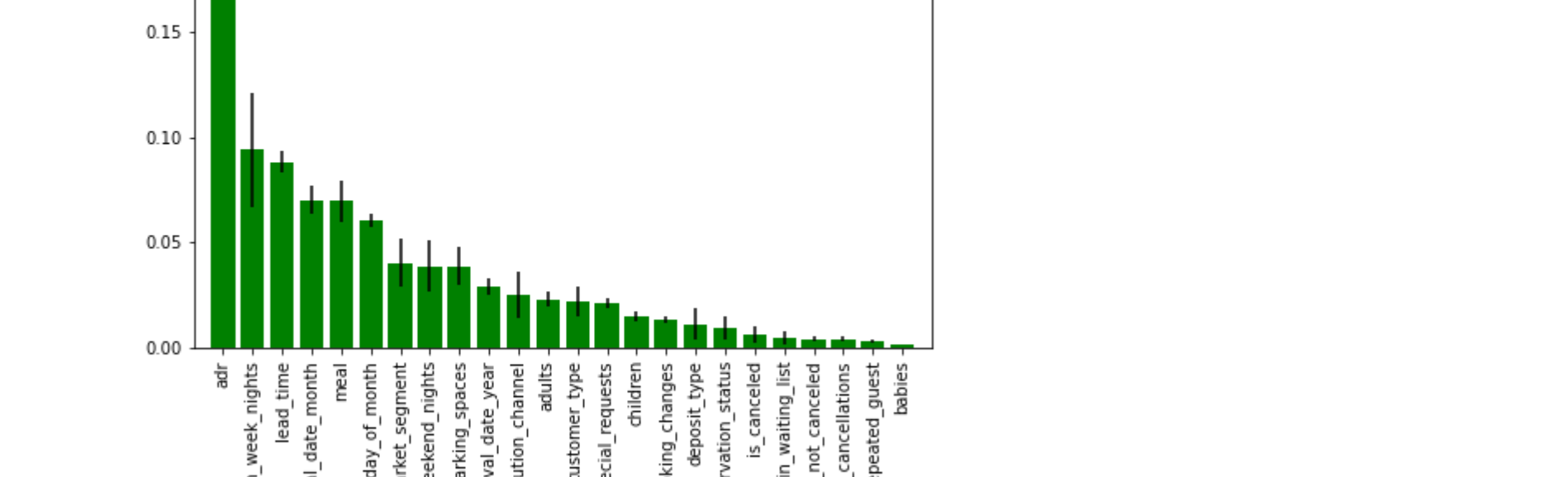



```
[43]: print("Feature ranking:")
for f in range(X_train.shape[1]):
    print("%d. feature %d (%2f) % (f + 1, indices[f]), importances(indices[f]))

Feature ranking:
1. feature 20 (0.305477)
2. feature 6 (0.093916)
3. feature 1 (0.088147)
4. feature 3 (0.070110)
5. feature 10 (0.069152)
6. feature 4 (0.060717)
7. feature 11 (0.040345)
8. feature 5 (0.038721)
9. feature 21 (0.038452)
10. feature 2 (0.029089)
11. feature 12 (0.025174)
12. feature 7 (0.023088)
13. feature 19 (0.021813)
14. feature 22 (0.021399)
15. feature 8 (0.014975)
16. feature 16 (0.013429)
17. feature 17 (0.011353)
18. feature 23 (0.009392)
19. feature 0 (0.006285)
20. feature 18 (0.004780)
21. feature 15 (0.004295)
22. feature 14 (0.004279)
23. feature 13 (0.003171)
24. feature 9 (0.001630)
```



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 assign 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 [45]:

```
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 [46]:

```
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',
'arrival_date_day_of_month', 'stays_in_week_nights',
'adults', 'meal', 'market_segment', 'ad',
'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 [47]:

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

In [48]:

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

In [49]:

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

In [50]:

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

In [51]:

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

In [52]:

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

Out [52]: 0.7288921152557516

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

In [53]:

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

Out [53]:

In [54]:

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

TP: 3893
FN: 8028
FP: 1682
TN: 22213

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 [55]:

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

Out [55]:

hotel	lead_time	arrival_date_year	arrival_date_month	arrival_date_day_of_month	stays_in_weekend_nights	stays_in_week_nights	adults
0	109.741108	2016.17344	6.566400	15.787094	0.706187	2.182954	1.8501
1	92.675686	2016.121443	6.544583	15.821243	1.189815	3.128732	1.8607

In conclusion, in both tests, the accuracy of the algorithm was around 0.72. The feature that has the most clarity about the type of the hotel is the "ad" which is the daily rate, the "ad" 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 n 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 [56]:

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

Firstly, I need to find out how many clusters do 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 what 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 [57]:

```
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 [58]:

```
number_of_clusters(hotel_unsupervised)
```



As we can see, in the first plot, the sharper elbow break is when I divide the data into 6 clusters. However, when we look at the second plot, the highest silhouette is when I divide the data into 10 clusters although the silhouette score in each number of clusters is not high. Therefore, I will divide the data into 6 clusters.

Below are two functions:

- create_kmeans_classifier - creating KMeans algorithm.
- clusters_information - display each cluster and his statistic details.

In [59]:

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

In [60]:

```
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 [61]:

```
kmeans = create_kmeans_classifier(6)
kmeans.fit(hotel_unsupervised)
silhouette_score(hotel_unsupervised, kmeans.labels_)
```

Out [61]: 0.2557500178616348


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

0
is_canceled lead_time arrival_date_year arrival_date_month \
8 1.0 0.115322 0.0 0.545455
9 1.0 0.181764 0.0 0.545455
10 1.0 0.031208 0.0 0.545455
27 1.0 0.081411 0.0 0.545455
34 1.0 0.061038 0.0 0.545455
... ..
110276 1.0 0.179104 1.0 0.272727
111351 1.0 0.054277 0.0 0.452443
111920 1.0 0.009498 0.0 0.000000
111921 1.0 0.008141 1.0 0.545455
117291 1.0 0.000000 1.0 0.636364

arrival_date_day_of_month stays_in_weekend_nights \
8 0.000000 0.000000
9 0.000000 0.000000
10 0.000000 0.000000
27 0.000000 0.105263
34 0.033333 0.052632
... ..
110276 0.800000 0.000000
111351 0.133333 0.052632
111920 0.0 0.000000
111921 0.533333 0.052632
117291 0.033333 0.000000

stays_in_week_nights adults children babies ... \
8 0.06 0.036364 0.0 0.0 ...
9 0.04 0.036364 0.0 0.0 ...
10 0.04 0.036364 0.0 0.0 ...
27 0.10 0.036364 0.0 0.0 ...
34 0.06 0.054545 0.0 0.0 ...
... ..
110276 0.00 0.000000 0.0 0.0 ...
111351 0.00 0.018182 0.0 0.0 ...
111920 0.02 0.018182 0.0 0.0 ...
111921 0.04 0.018182 0.0 0.0 ...
117291 0.02 0.018182 0.0 0.0 ...

deposit_type_Non Refund deposit_type_Refundable \
8 0.0 0.0
9 0.0 0.0
10 0.0 0.0
27 0.0 0.0
34 0.0 0.0
... ..
110276 ... ...
111351 ... ...
111920 0.0 0.0
111921 0.0 0.0
117291 0.0 0.0

customer_type_Contract customer_type_Group customer_type_Transient \
8 0.0 0.0 0.0
9 0.0 0.0 0.0
10 0.0 0.0 1.0
11 0.0 0.0 1.0
27 0.0 0.0 0.0
34 0.0 0.0 1.0
... ..
110276 ... ...
111351 ... ...
111920 0.0 0.0
111921 0.0 0.0
117291 0.0 0.0

reservation_status_Check-Out reservation_status_No-Show label
8 0.0 0.0 0
9 0.0 0.0 0
10 0.0 0.0 0
27 0.0 0.0 0
34 0.0 0.0 0
... ..
110276 0.0 0.0 0
111351 0.0 0.0 0
111920 0.0 0.0 0
111921 0.0 0.0 1
117291 0.0 0.0 0

[27272 rows x 48 columns]
is_canceled lead_time arrival_date_year arrival_date_month \
count 22722.000000 22722.000000 22722.000000 22722.000000
mean 1.0 0.143573 0.671926 0.499516
std 0.0 0.114967 0.321782 0.264966
min 0.0 0.031764 0.000000 0.000000
25% 1.0 0.047490 0.000000 0.000000
50% 1.0 0.113976 0.500000 0.545455
75% 1.0 0.218433 1.000000 0.636364
max 1.0 0.590232 1.000000 1.000000

arrival_date_day_of_month stays_in_week_nights \
count 22722.000000 22722.000000
mean 0.496811 0.058610
std 0.296879 0.055687
min 0.000000 0.000000
25% 0.233333 0.000000
50% 0.500000 0.052632
75% 0.766667 0.105263
max 1.000000 0.842105

stays_in_week_nights adults children babies ... \
count 22722.000000 22722.000000 22722.000000 22722.000000
mean 0.058273 0.035953 0.018731 0.000572
std 0.041962 0.010631 0.052902 0.007601
min 0.000000 0.000000 0.000000 0.000000
25% 0.040000 0.036364 0.000000 0.000000
50% 0.060000 0.036364 0.000000 0.000000
75% 0.080000 0.036364 0.000000 0.000000
max 0.800000 0.472727 0.300000 0.200000

deposit_type_Non Refund deposit_type_Refundable \
count 22722.000000 22722.000000
mean 0.001232 0.000792
std 0.030883 0.028135
min 0.000000 0.000000
25% 0.000000 0.000000
50% 0.000000 0.000000
75% 0.000000 0.000000
max 1.000000 1.000000

customer_type_Contract customer_type_Group customer_type_Transient \
count 22722.000000 22722.000000 22722.000000
mean 0.023342 0.001672 0.000000
std 0.152870 0.040862 0.204632
min 0.000000 0.000000 0.000000
25% 0.000000 0.036364 0.000000
50% 0.000000 0.000000 1.000000
75% 0.000000 0.036364 1.000000
max 1.000000 1.000000 1.000000

customer_type_Transient-Party reservation_status_Canceled \
count 22722.000000 22722.000000
mean 0.133591 0.173037
std 0.335917 0.473037
min 0.000000 0.000000
25% 0.000000 0.000000
50% 0.000000 1.000000
75% 0.000000 1.000000
max 1.000000 1.000000

reservation_status_Check-Out reservation_status_No-Show label
count 22722.000000 22722.000000
mean 0.0 0.0
std 0.0 0.0
min 0.0 0.0
25% 0.0 0.0
50% 0.0 0.0
75% 0.0 0.0
max 0.0 0.0

[8 rows x 48 columns]
is_canceled lead_time arrival_date_year arrival_date_month \
3 0.0 0.017639 0.0 0.545455
4 0.0 0.018996 0.0 0.545455
9 0.0 0.018996 0.0 0.545455
11 0.0 0.047490 0.0 0.545455
12 0.0 0.092266 0.0 0.545455
... ..
40052 0.280667 0.000000 0.000000
40053 0.364993 1.0 0.636364
40055 0.287653 1.0 0.636364
40058 0.286296 1.0 0.636364
40059 0.218453 1.0 0.636364

arrival_date_day_of_month stays_in_weekend_nights \
3 0.000000 0.000000
4 0.000000 0.000000
9 0.000000 0.000000
11 0.000000 0.000000
12 0.000000 0.000000
... ..
40052 0.000000 0.000000
40053 0.000000 0.000000
40055 1.000000 0.105263
40058 1.000000 0.215262
40059 1.000000 0.842105

stays_in_week_nights adults children babies ... \
3 0.0 0.036364 0.0 0.0 ...
4 0.04 0.036364 0.0 0.0 ...
9 0.04 0.036364 0.0 0.0 ...
11 0.08 0.036364 0.0 0.0 ...
12 0.08 0.036364 0.0 0.0 ...
... ..
40052 0.20 0.036364 0.0 0.0 ...
40053 0.20 0.036364 0.0 0.0 ...
40055 0.16 0.036364 0.1 0.0 ...
40058 0.20 0.036364 0.0 0.0 ...
40059 0.20 0.036364 0.0 0.0 ...

deposit_type_Non Refund deposit_type_Refundable \
3 0.0 0.0
4 0.0 0.0
9 0.0 0.0
11 0.0 0.0
12 0.0 0.0
... ..
40052 0.000000 0.000000
40053 0.000000 0.000000
40055 0.000000 0.000000
40058 0.000000 0.000000
40059 0.000000 0.000000

customer_type_Contract customer_type_Group customer_type_Transient \
3 0.0 0.0 1.0
4 0.0 0.0 1.0
9 0.0 0.0 1.0
11 0.0 0.0 1.0
12 0.0 0.0 1.0
... ..
40052 0.0 0.0 0.0
40053 1.0 0.0 0.0
40055 0.0 0.0 0.0
40058 1.0 0.0 0.0
40059 0.0 0.0 0.0

customer_type_Transient-Party reservation_status_Canceled \
3 0.0 0.0
4 0.0 0.0
9 0.0 0.0
11 0.0 0.0
12 0.0 0.0
... ..
40052 0.0 0.0
40053 0.0 0.0
40055 0.0 0.0
40058 0.0 0.0
40059 0.0 0.0

reservation_status_Check-Out reservation_status_No-Show label
3 1.0 0.0 1
4 1.0 0.0 1
9 0.0 0.0 0
11 1.0 0.0 1
12 1.0 0.0 1
... ..
40052 1.0 0.0 1
40053 1.0 0.0 1
40055 1.0 0.0 1
40058 1.0 0.0 1
40059 1.0 0.0 1

[18192 rows x 48 columns]
is_canceled lead_time arrival_date_year arrival_date_month \
count 18792.000000 18792.000000 18792.000000 18792.000000
mean 0.003735 0.102016 0.557125 0.501914
std 0.005021 0.117589 0.001788 0.210462
min 0.000000 0.000000 0.000000 0.000000
25% 0.000000 0.000000 0.000000 0.000000
50% 0.000000 0.000000 0.000000 0.000000
75% 0.000000 0.000000 0.000000 0.000000
max 1.000000 0.166893 1.000000 0.727273

arrival_date_day_of_month stays_in_weekend_nights \
count 18792.000000 18792.000000
mean 0.493668 0.056100
std 0.295468 0.061729
min 0.000000 0.000000
25% 0.233333 0.000000
50% 0.500000 0.052632
75% 0.766667 0.105263
max 1.000000 0.842105

stays_in_week_nights adults children babies ... \
count 18792.000000 18792.000000 18792.000000 18792.000000
mean 0.063841 0.033883 0.011383 0.001453
min 0.000000 0.000000 0.000000 0.000000
25% 0.020000 0.036364 0.000000 0.000000
50% 0.060000 0.036364 0.000000 0.000000
75% 0.100000 0.036364 0.000000 0.000000
max 0.800000 0.072727 1.000000 0.200000

deposit_type_Non Refund deposit_type_Refundable \
count 18792.000000 18792.000000
mean 0.000053 0.000426
std 0.007295 0.000029
min 0.000000 0.000000
25% 0.000000 0.000000
50% 0.000000 0.000000
75% 0.000000 0.000000
max 1.000000 1.000000

customer_type_Contract customer_type_Group customer_type_Transient \
count 18792.000000 18792.000000 18792.000000
mean 0.085834 0.010309 0.834291
std 0.280127 0.103871 0.000000
min 0.000000 0.000000 0.000000
25% 0.000000 0.000000 1.000000
50% 0.000000 0.000000 1.000000
75% 0.000000 0.000000 1.000000
max 1.000000 1.000000 1.000000

customer_type_Transient-Party reservation_status_Canceled \
count 18792.000000 18792.000000
mean 0.068966 0.000000
std 0.351176 0.000000
min 0.000000 0.000000
25% 0.000000 0.000000
50% 0.000000 0.000000
75% 0.000000 0.000000
max 1.000000 0.000000

reservation_status_Check-Out reservation_status_No-Show label
count 18792.000000 18792.000000
mean 0.096275 0.003725
std 0.233321 0.000000
min 0.000000 0.000000
25% 1.000000 0.000000
50% 1.000000 0.000000
75% 1.000000 0.000000
max 1.000000 1.000000

[18192 rows x 48 columns]
is_canceled lead_time arrival_date_year arrival_date_month \
count 21331.000000 21331.000000 21331.000000 21331.000000
mean 0.224040 0.186384 0.461394 0.252298
std 0.416958 0.163469 0.361954 0.282928
min 0.000000 0.000000 0.000000 0.000000
25% 0.000000 0.050204 0.000000 0.272727
50% 0.000000 0.138399 0.500000 0.545455
75% 0.000000 0.280988 0.500000 0.636364
max 1.000000 0.735414 1.000000 1.000000

arrival_date_day_of_month stays_in_weekend_nights \
count 21331.000000 21331.000000
mean 0.498944 0.042352
std 0.284473 0.048352
min 0.000000 0.000000
25% 0.266667 0.000000
50% 0.500000 0.052632
75% 0.733333 0.105263
max 1.000000 0.842105

stays_in_week_nights adults children babies ... \
count 21331.000000 21331.000000 21331.000000 21331.000000
mean 0.044954 0.031076 0.001814 0.000028
std 0.032545 0.009231 0.016433 0.007742
min 0.000000 0.000000 0.000000 0.000000
25% 0.020000 0.036364 0.000000 0.000000
50% 0.040000 0.036364 0.000000 0.000000
75% 0.060000 0.036364 0.000000 0.000000
max 0.700000 0.072727 1.000000 0.500000

deposit_type_Non Refund deposit_type_Refundable \
count 21331.000000 21331.000000
mean 0.004172 0.005954
std 0.064460 0.076932
min 0.000000 0.000000
25% 0.000000 0.000000
50% 0.000000 0.000000
75% 0.000000 0.000000
max 1.000000 1.000000

customer_type_Contract customer_type_Group customer_type_Transient \
count 21331.000000 21331.000000 21331.000000
mean 0.004554 0.005626 0.000000
std 0.066588 0.074795 0.000000
min 0.000000 0.000000 0.000000
25% 0.000000 0.000000 0.000000
50% 0.000000 0.000000 0.000000
75% 0.000000 0.000000 0.000000
max 1.000000 1.000000 1.000000

customer_type_Transient-Party reservation_status_Canceled \
count 21331.000000 21331.000000
mean 0.098921 0.214992
std 0.099890 0.410827
min 0.000000 0.000000
25% 0.000000 0.000000
50% 0.000000 0.000000
75% 1.000000 0.000000
max 1.000000 1.000000

reservation_status_Check-Out reservation_status_No-Show label
count 21331.000000 21331.000000
mean 0.041698 0.094691
std 0.100000 0.000000
min 0.000000 0.000000
25% 1.000000 0.000000
50% 1.000000 0.000000
75% 1.000000 0.000000
max 1.000000 2.000000

[8 rows x 48 columns]
is_canceled lead_time arrival_date_year arrival_date_month \
count 40060.000000 40060.000000 40060.000000 40060.000000
mean 0.0 0.081411 0.0 0.545455
std 0.0 0.081411 0.0 0.545455
min 0.0 0.081411 0.0 0.545455
25% 0.0 0.081411 0.0 0.545455
50% 0.0 0.081411 0.0 0.545455
75% 0.0 0.081411 0.0 0.545455
max 0.0 0.081411 0.0 0.545455

arrival_date_day_of_month stays_in_weekend_nights \
40060 0.0 0.000000 0.000000
40085 0.066667 0.000000
40092 0.000000 0.000000
40113 0.200000 0.105263
40119 0.233333 0.052632
... ..
119381 0.000000 0.000000
119382 0.000000 0.000000
119383 0.000000 0.000000
119384 0.000000 0.000000
119385 0.000000 0.000000

stays_in_week_nights adults children babies ... \
40060 0.0 0.036364 0.0 0.0 ...
40085 0.10 0.036364 0.0 0.0 ...
40104 0.04 0.036364 0.1 0.0 ...
40113 0.16 0.036364 0.0 0.0 ...
40119 0.08 0.036364 0.0 0.0 ...
... ..
119381 0.0 0.036364 0.0 0.0 ...
119382 0.10 0.054545 0.0 0.0 ...
119383 0.10 0.036364 0.0 0.0 ...
119384 0.10 0.036364 0.0 0.0 ...
119385 0.10 0.036364 0.0 0.0 ...

deposit_type_Non Refund deposit_type_Refundable \
40060 0.0 0.0
40085 0.0 0.0
40104 0.0 0.0
40113 0.0 0.0
40119 0.0 0.0
... ..
119381 0.0 0.0
119382 0.0 0.0
119383 0.0 0.0
119384 0.0 0.0
119385 0.0 0.0

customer_type_Contract customer_type_Group customer_type_Transient \
40060 0.0 0.0 1.0
40085 0.0 0.0 1.0
40104 0.0 0.0 1.0
40113 0.0 0.0 1.0
40119 0.0 0.0 1.0
... ..
119381 0.0 0.0 1.0
119382 0.0 0.0 1.0
119383 0.0 0.0 1.0
119384 0.0 0.0 1.0
119385 0.0 0.0 1.0

customer_type_Transient-Party reservation_status_Canceled \
40060 0.0 0.0
40085 0.0 0.0
40104 0.0 0.0
40113 0.0 0.0
40119 0.0 0.0
... ..
119381 0.0 0.0
119382 0.0 0.0
119383 0.0 0.0
119384 0.0 0.0
119385 0.0 0.0

reservation_status_Check-Out reservation_status_No-Show label
count 40060.000000 40060.000000
mean 0.0 0.0
std 0.0 0.0
min 0.0 0.0
25% 0.0 0.0
50% 0.0 0.0
75% 0.0 0.0
max 0.0 0.0

[8 rows x 48 columns]
is_canceled lead_time arrival_date_year arrival_date_month \
count 14703.000000 14703.000000 14703.000000 14703.000000
mean 0.0 0.181463 0.0 0.636364
std 0.0 0.093541 0.321347 0.281916
min 0.0 0.000000 0.000000 0.000000
25% 0.0 0.019452 0.000000 0.272727
50% 0.0 0.036988 0.000000 0.545455
75% 0.0 0.138399 0.000000 0.727273
max 0.0 0.582090 1.000000 1.000000

arrival_date_day_of_month stays_in_weekend_nights \
count 28723.000000 28723.000000
mean 0.489402 0.032507
std 0.294238 0.045759
min 0.000000 0.000000
25% 0.233333 0.000000
50% 0.500000 0.052632
75% 0.733333 0.105263
max 1.000000 0.842105

stays_in_week_nights adults children babies ... \
count 28723.000000 28723.000000 28723.000000 28723.000000
mean 0.041802 0.030001 0.000027 0.000000
std 0.022936 0.007113 0.002020 0.000000
min 0.000000 0.000000 0.000000 0.000000
25% 0.020000 0.036364 0.000000 0.000000
50% 0.040000 0.036364 0.000000 0.000000
75% 0.060000 0.036364 0.000000 0.000000
max 0.820000 0.072727 0.300000 0.000000

deposit_type_Non Refund deposit_type_Refundable \
count 28723.000000 28723.000000
mean 0.009881 0.000408
std 0.125938 0.020198
min 0.000000 0.000000
25% 0.000000 0.000000
50% 1.000000 0.000000
75% 1.000000 0.000000
max 1.000000 1.000000

customer_type_Contract customer_type_Group customer_type_Transient \
count 28723.000000 28723.000000 28723.000000
mean 0.038332 0.003515 0.953173
std 0.000000 0.000000 0.000000
min 0.000000 0.000000 0.000000
25% 0.000000 0.000000 1.000000
50% 0.000000 0.000000 1.000000
75% 0.000000 0.000000 1.000000
max 1.000000 1.000000 1.000000

customer_type_Transient-Party reservation_status_Canceled \
count 28723.000000 28723.000000
mean 0.004979 0.000000
std 0.077184 0.000000
min 0.000000 0.000000
25% 0.000000 0.000000
50% 0.000000 0.000000
75% 0.000000 0.000000
max 1.000000 0.000000

reservation_status_Check-Out reservation_status_No-Show label
count 28723.000000 28723.000000
mean 0.0 0.0
std 0.0 0.0
min 0.0 0.0
25% 0.0 0.0
50% 0.0 0.0
75% 0.0 0.0
max 0.0 0.0

[14703 rows x 48 columns]
is_canceled lead_time arrival_date_year arrival_date_month \
count 14703.000000 14703.000000 14703.000000 14703.000000
mean 0.0 0.290561 0.508600 0.512419
std 0.0 0.137993 0.368590 0.276918
min 0.0 0.000000 0.000000 0.000000
25% 1.0 0.141113 0.000000 0.272727
50% 1.0 0.252374 0.500000 0.454545
75% 1.0 0.412453 1.000000 0.636364
max 1.0 0.853460 1.000000 1.000000

arrival_date_day_of_month stays_in_weekend_nights \
count 14703.000000 14703.000000
mean 0.486975 0.032507
std 0.289724 0.043046
min 0.000000 0.000000
25% 0.233333 0.000000
50% 0.500000 0.052632
75% 0.733333 0.105263
max 1.000000 0.368421

stays_in_week_nights adults children babies ... \
count 14703.000000 14703.000000 14703.000000 14703.000000
mean 0.041802 0.030001 0.000027 0.000000
std 0.022936 0.007113 0.002020 0.000000
min 0.000000 0.000000 0.000000 0.000000
25% 0.020000 0.036364 0.000000 0.000000
50% 0.040000 0.036364 0.000000 0.000000
75% 0.060000 0.036364 0.000000 0.000000
max 0.820000 0.072727 0.300000 0.000000

deposit_type_Non Refund deposit_type_Refundable \
count 14703.000000 14703.000000
mean 0.009881 0.000408
std 0.125938 0.020198
min 0.000000 0.000000
25% 0.000000 0.000000
50% 1.000000 0.000000
75% 1.000000 0.000000
max 1.000000 1.000000

customer_type_Contract customer_type_Group customer_type_Transient \
count 14703.000000 14703.000000 14703.000000
mean 0.048085 0.000000 0.880841
std 0.213954 0.000000 0.323987
min 0.000000 0.000000 0.000000
25% 0.000000 0.000000 1.000000
50% 0.000000 0.000000 1.000000
75% 0.000000 0.000000 1.000000
max 1.000000 0.000000 1.000000

customer_type_Transient-Party reservation_status_Canceled \
count 14703.000000 14703.000000
mean 0.071074 0.098844
std 0.259757 0.033985
min 0.000000 0.000000
25% 0.000000 0.000000
50% 0.000000 0.000000
75% 0.000000 0.000000
max 1.000000 0.000000

reservation_status_Check-Out reservation_status_No-Show label

```



```
0.966667 0.052632
119364 0.800000 0.157895
119365 1.000000 0.105263
119377

stays_in_week_nights adults children babies ... \
0 0.00 0.036364 0.0 0.0 ...
1 0.00 0.036364 0.0 0.0 ...
2 0.02 0.018182 0.0 0.0 ...
6 0.04 0.036364 0.0 0.0 ...
7 0.04 0.036364 0.0 0.0 ...
...
119354 0.08 0.036364 0.0 0.0 ...
119361 0.10 0.018182 0.0 0.0 ...
119364 0.08 0.036364 0.0 0.0 ...
119365 0.14 0.036364 0.0 0.0 ...
119377 0.06 0.036364 0.0 0.0 ...

deposit_type_Non Refund deposit_type_Refundable \
0 0.0 0.0
1 0.0 0.0
2 0.0 0.0
6 0.0 0.0
7 0.0 0.0
...
119354 0.0 0.0
119361 0.0 0.0
119364 0.0 0.0
119365 0.0 0.0
119377 0.0 0.0

customer_type_Contract customer_type_Group customer_type_Transient \
0 0.0 0.0 1.0
1 0.0 0.0 1.0
2 0.0 0.0 1.0
6 0.0 0.0 1.0
7 0.0 0.0 1.0
...
119354 0.0 0.0 1.0
119361 0.0 0.0 1.0
119364 0.0 0.0 1.0
119365 0.0 0.0 1.0
119377 0.0 0.0 1.0

customer_type_Transient-Party reservation_status_Canceled \
0 0.0 0.0
1 0.0 0.0
2 0.0 0.0
6 0.0 0.0
7 0.0 0.0
...
119354 0.0 0.0
119361 0.0 0.0
119364 0.0 0.0
119365 0.0 0.0
119377 0.0 0.0

reservation_status_Check-Out reservation_status_No-Show label
0 1.0 0.0 5
1 1.0 0.0 5
2 1.0 0.0 5
6 1.0 0.0 5
7 1.0 0.0 5
...
119354 1.0 0.0 ...
119361 1.0 0.0 5
119364 1.0 0.0 5
119365 1.0 0.0 5
119377 1.0 0.0 5

[13115 rows x 48 columns]
is_canceled lead_time arrival_date_year arrival_date_month \
count 13115.000000 13115.000000 13115.000000 13115.000000
mean 0.148335 0.965632 0.583568 0.495228
std 0.35549 0.093518 0.360882 0.292798
min 0.00000 0.000000 0.000000 0.000000
25% 0.00000 0.001357 0.500000 0.272727
50% 0.00000 0.018896 0.500000 0.545455
75% 0.00000 0.097693 1.000000 0.727273
max 1.00000 1.000000 1.000000 1.000000

arrival_date_day_of_month stays_in_weekend_nights \
count 13115.000000 13115.000000
mean 0.497355 0.043670
std 0.295283 0.053497
min 0.000000 0.000000
25% 0.233333 0.000000
50% 0.500000 0.052632
75% 0.766667 0.105263
max 1.000000 1.000000

stays_in_week_nights adults children babies ... \
count 13115.000000 13115.000000 13115.000000 13115.000000 ...
mean 0.041191 0.037438 0.050192 0.013565 ...
std 0.041191 0.037438 0.050192 0.013565 ...
min 0.000000 0.000000 0.000000 0.000000 ...
25% 0.020000 0.036364 0.000000 0.000000 ...
50% 0.040000 0.036364 0.000000 0.000000 ...
75% 0.060000 0.036364 0.000000 0.000000 ...
max 1.000000 1.000000 0.300000 0.200000 ...

deposit_type_Non Refund deposit_type_Refundable \
count 13115.000000 13115.000000
mean 0.000229 0.000229
std 0.015123 0.015123
min 0.000000 0.000000
25% 0.000000 0.000000
50% 0.000000 0.000000
75% 0.000000 0.000000
max 1.000000 1.000000

customer_type_Contract customer_type_Group customer_type_Transient \
count 13115.000000 13115.000000 13115.000000
mean 0.001220 0.008616 0.905757
std 0.034908 0.092426 0.292178
min 0.000000 0.000000 0.000000
25% 0.000000 0.000000 0.000000
50% 0.000000 0.000000 1.000000
75% 0.000000 0.000000 1.000000
max 1.000000 1.000000 1.000000

customer_type_Transient-Party reservation_status_Canceled \
count 13115.000000 13115.000000
mean 0.084407 0.131224
std 0.278008 0.337658
min 0.000000 0.000000
25% 0.000000 0.000000
50% 0.000000 0.000000
75% 0.000000 0.000000
max 1.000000 1.000000

reservation_status_Check-Out reservation_status_No-Show label
count 13115.000000 13115.000000 13115.0
mean 0.85162 0.017156 5.0
std 0.35549 0.129857 0.0
min 0.00000 0.000000 5.0
25% 1.00000 0.000000 5.0
50% 1.00000 0.000000 5.0
75% 1.00000 0.000000 5.0
max 1.00000 1.000000 5.0

[8 rows x 48 columns]
```

As we can see [here](#), the silhouette score that we got is very low. An explanation for this is because the values of the features are very similar that the algorithm had difficulties to handle those samples and to separate them into groups that have the same behavior.