

Week 7 Assignment

DSC 630

Bellevue University

Yograj Karki

DSC630, Assignment 7.2: Create Optimal Hotel Recommendations

Online travel agencies are scrambling to meet the artificial intelligence driven personalization standard set by companies like Amazon and Netflix. In addition, the world of online travel has become a highly competitive space where brands try to capture user attention (and wallet) with recommendations, comparing, matching, and sharing. For this assignment, we would like to create the optimal hotel recommendations for Expedia's users that are searching for a hotel to book. For this assignment, you need to predict which "hotel cluster" the user is likely to book, given his (or her) search details.

The data set can be found at Kaggle: Expedia Hotel Recommendations. To get started, I would suggest exploring the file train.csv, which contains the logs of user behavior. There is another file named destination.csv, which contains information related to hotel reviews made by users. There is a lot of data here, and making an accurate prediction is rather difficult, e.g., simply running a standard prediction algorithm will probably yield below 10% accuracy. Start by doing some exploratory analysis of this data to help understand how to make a prediction on the hotel cluster the user is likely to select. Then, split train.csv into a training and test set (feel free to select a smaller random subset of train.csv). Then, build at least two prediction models from the training set, and report the accuracies on the test set. As I mentioned, this is a difficult problem, so be creative with your solutions. You might want to try building your own predictor rather than a standard predictor model, e.g., a random forest. The purpose of this project is not necessarily to get great results but to understand the nuances and challenges of such problems.

Exploratory Data Analysis

```
In [1]: # Loading Libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

In [2]: # read CSV file
# filename = "expedia-hotel-recommendations/train.csv"
df = pd.read_csv(filename)

# count no. of lines
df.info()
```

Sampling the rows of the csv file since it's very large (taking 6.7 GB+ memory)

```
In [3]: import random
random.seed(7)
import random

n = 37670293 # number of records in file
s = 100000 # arbitrary sample size
filename = "expedia-hotel-recommendations/train.csv"
skip = int(len(random.sample(range(1), n), s)) # skip rows
df = pd.read_csv(filename, skiprows=skip)

In [4]: df.shape
Out[4]: (100000, 24)

In [5]: df.info()
Out[5]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 100000 entries, 0 to 99999
Data columns (total 24 columns):
#   Column                Non-Null Count  Dtype
---  --
0   date_time              100000 non-null object
1   site_name              100000 non-null int64
2   posa_continent         100000 non-null int64
3   user_location_country  100000 non-null int64
4   user_location_region   100000 non-null int64
5   user_location_city     100000 non-null int64
6   orig_destination_distance 64037 non-null float64
7   user_id                100000 non-null int64
8   is_mobile              100000 non-null int64
9   is_package             100000 non-null int64
10  arch_channel           100000 non-null object
11  arch_ci                99999 non-null object
12  arch_co                100000 non-null object
13  arch_children_cnt      100000 non-null int64
14  arch_children_cnt      100000 non-null int64
15  arch_rm_cnt            100000 non-null int64
16  arch_destination_id    100000 non-null int64
17  is_booking             100000 non-null int64
18  hotel_continent        100000 non-null int64
19  cnt                    100000 non-null int64
20  hotel_continent        100000 non-null int64
21  hotel_country          100000 non-null int64
22  hotel_market           100000 non-null int64
23  hotel_cluster          100000 non-null int64
dtypes: float64(1), int64(20), object(3)
memory usage: 18.3+ MB

In [6]: df.head()
Out[6]:
```

	date_time	site_name	posa_continent	user_location_country	user_location_region	user_location_city	orig_destination_distance	user_id	is_
0	2014-05-17 17:42:23	23	30	4	195	991	47725	NaN	1048
1	2014-05-12 22:20:03	2	3	66	174	21356	2103.8393	7523	
2	2013-04-02 19:50:30	24	2	3	50	5703	NaN	9616	
3	2014-08-12 14:54:54	2	3	66	348	48862	3864.2730	10706	
4	2014-08-26 07:53:13	2	3	66	322	48585	224.9042	14864	

5 rows x 24 columns

```
In [7]: # Drop any rows containing null values
df.dropna(inplace=True)

In [8]: df.info()
Out[8]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 64037 entries, 1 to 99999
Data columns (total 24 columns):
#   Column                Non-Null Count  Dtype
---  --
0   date_time              64037 non-null object
1   site_name              64037 non-null int64
2   posa_continent         64037 non-null int64
3   user_location_country  64037 non-null int64
4   user_location_region   64037 non-null int64
5   user_location_city     64037 non-null int64
6   orig_destination_distance 64037 non-null float64
7   user_id                64037 non-null int64
8   is_mobile              64037 non-null int64
9   is_package             64037 non-null int64
10  arch_channel           64037 non-null object
11  arch_ci                64037 non-null object
12  arch_co                64037 non-null object
13  arch_children_cnt      64037 non-null int64
14  arch_children_cnt      64037 non-null int64
15  arch_rm_cnt            64037 non-null int64
16  arch_destination_id    64037 non-null int64
17  is_booking             64037 non-null int64
18  hotel_continent        64037 non-null int64
19  cnt                    64037 non-null int64
20  hotel_continent        64037 non-null int64
21  hotel_country          64037 non-null int64
22  hotel_market           64037 non-null int64
23  hotel_cluster          64037 non-null int64
dtypes: float64(1), int64(20), object(3)
memory usage: 12.2+ MB

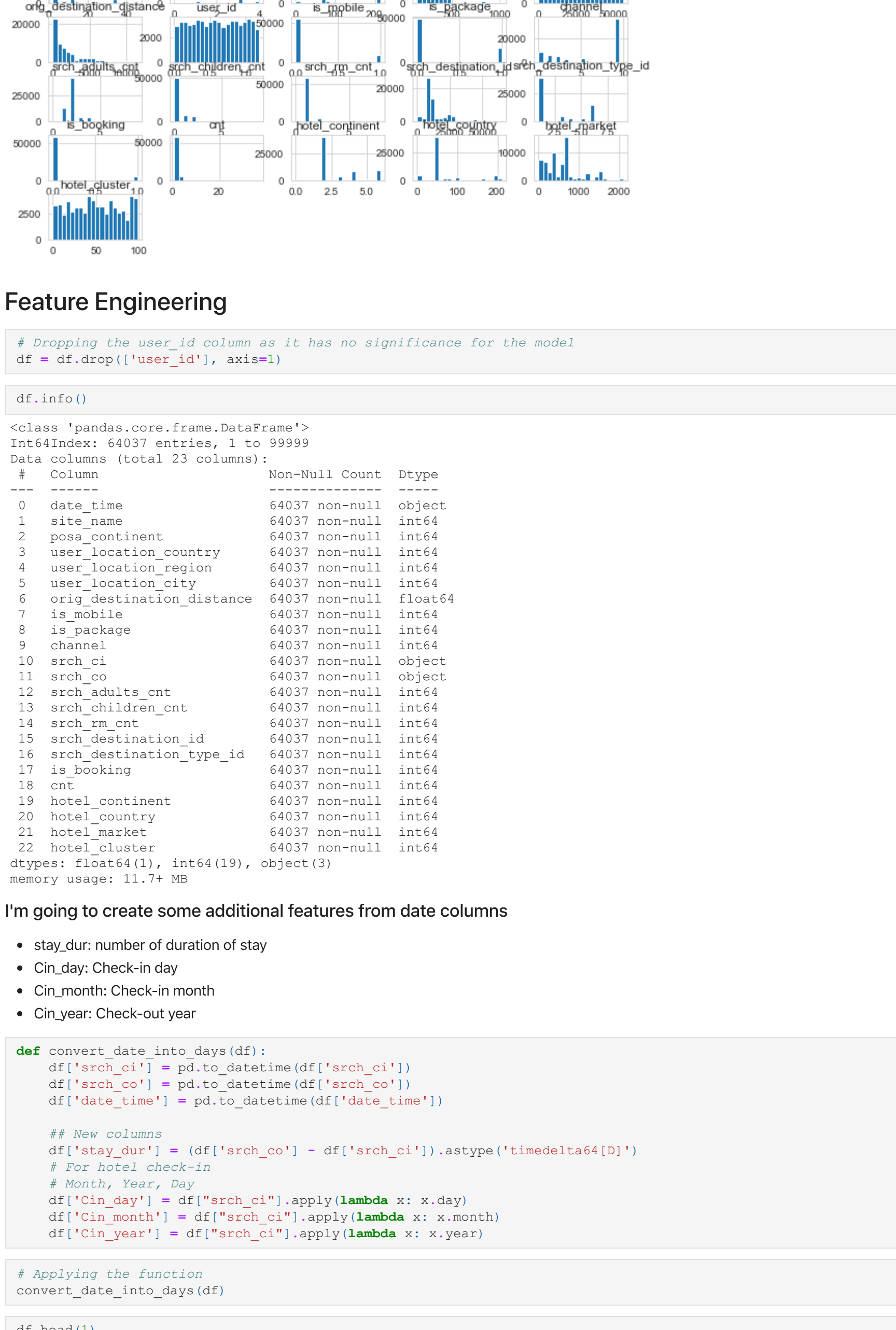
In [9]: # Load Libraries
import matplotlib.pyplot as plt
import seaborn as sns

In [10]: # Set up chart style
sns.set_style("whitegrid")

# Set up the figure size
plt.rcParams["figure.figsize"] = (10, 5)

In [11]: # Plot frequency for each hotel_clusters
df["hotel_cluster"].value_counts().plot(kind='bar', colormap="Set3", figsize=(15, 5))

Out[11]: <AxesSubplot>
```

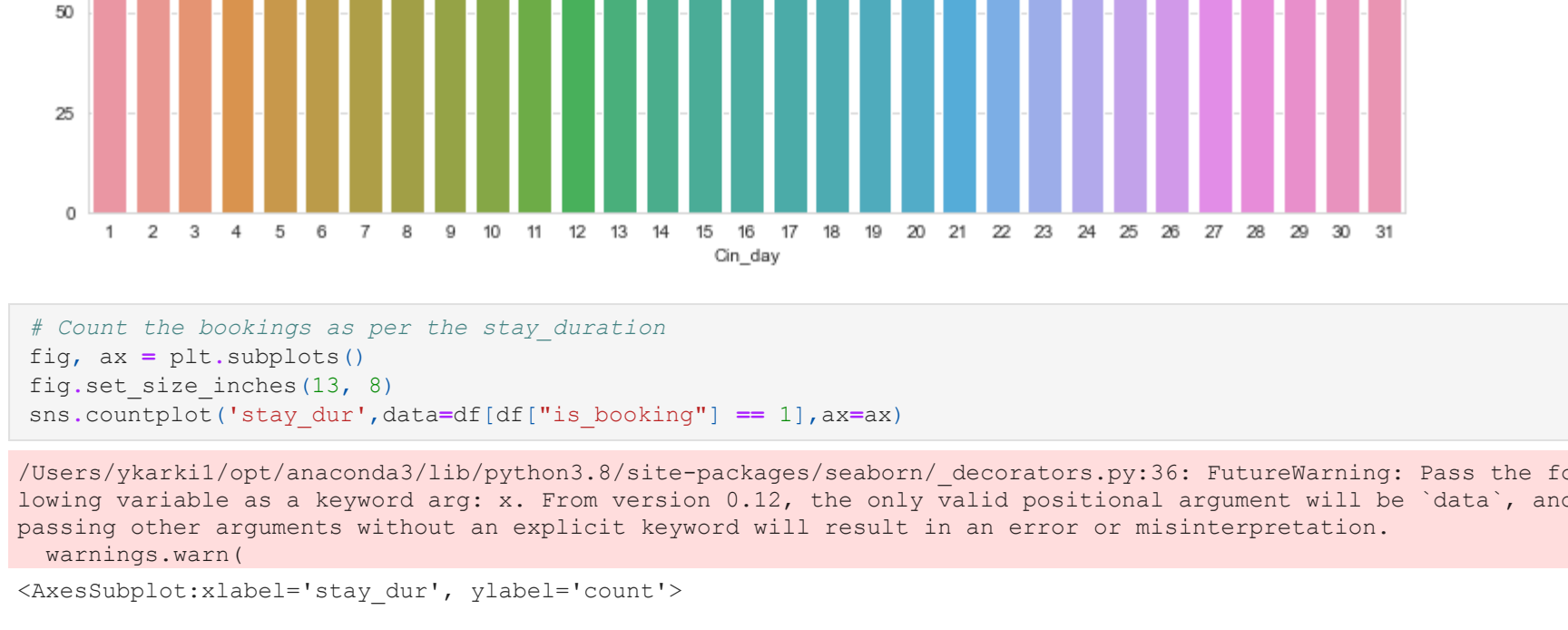


```
In [12]: # heatmap
fig, ax = plt.subplots()
fig.set_size_inches(15, 10)
sns.heatmap(df.corr(), cmap="coolwarm", ax=ax, annot=True, linewidths=2)

Out[12]: <AxesSubplot>
```



```
In [13]: # plot histograms
hist = df.hist(bins=20)
```



Feature Engineering

```
In [14]: # Dropping the user_id column as it has no significance for the model
df = df.drop(["user_id"], axis=1)

In [15]: df.info()
Out[15]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 64037 entries, 1 to 99999
Data columns (total 23 columns):
#   Column                Non-Null Count  Dtype
---  --
0   date_time              64037 non-null object
1   site_name              64037 non-null int64
2   posa_continent         64037 non-null int64
3   user_location_country  64037 non-null int64
4   user_location_region   64037 non-null int64
5   user_location_city     64037 non-null int64
6   orig_destination_distance 64037 non-null float64
7   user_id                64037 non-null int64
8   is_mobile              64037 non-null int64
9   is_package             64037 non-null int64
10  arch_channel           64037 non-null object
11  arch_ci                64037 non-null object
12  arch_co                64037 non-null object
13  arch_children_cnt      64037 non-null int64
14  arch_rm_cnt            64037 non-null int64
15  arch_destination_id    64037 non-null int64
16  arch_destination_type_id 64037 non-null int64
17  is_booking             64037 non-null int64
18  cnt                    64037 non-null int64
19  hotel_continent        64037 non-null int64
20  hotel_country          64037 non-null int64
21  hotel_market           64037 non-null int64
22  hotel_cluster          64037 non-null int64
dtypes: float64(1), int64(19), object(3)
memory usage: 11.7+ MB
```

I'm going to create some additional features from data columns

- stay_dur: number of duration of stay
- Cin_day: Check-in day
- Cin_month: Check-in month
- Cin_year: Check-out year

```
In [16]: def convert_date_into_days(df):
df['arch_ci'] = pd.to_datetime(df['arch_ci'])
df['arch_co'] = pd.to_datetime(df['arch_co'])
df['date_time'] = pd.to_datetime(df['date_time'])

# New columns
df['stay_dur'] = df['arch_co'] - df['arch_ci'].astype('timedelta64[D]')
df['Cin_day'] = df['Cin_day']
df['Cin_month'] = df['Cin_month']
df['Cin_year'] = df['Cin_year']

In [17]: # Applying the function
convert_date_into_days(df)

In [18]: df.head(1)
```

```
Out[18]:
```

	date_time	site_name	posa_continent	user_location_country	user_location_region	user_location_city	orig_destination_distance	is_mobile	is_package
1	2014-05-12 22:20:03	2	3	66	174	21356	2103.8393	1	

1 rows x 27 columns

```
In [19]: df.info()
Out[19]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 64037 entries, 1 to 99999
Data columns (total 27 columns):
#   Column                Non-Null Count  Dtype
---  --
0   date_time              64037 non-null object
1   site_name              64037 non-null int64
2   posa_continent         64037 non-null int64
3   user_location_country  64037 non-null int64
4   user_location_region   64037 non-null int64
5   user_location_city     64037 non-null int64
6   orig_destination_distance 64037 non-null float64
7   user_id                64037 non-null int64
8   is_mobile              64037 non-null int64
9   is_package             64037 non-null int64
10  arch_channel           64037 non-null object
11  arch_ci                64037 non-null object
12  arch_co                64037 non-null object
13  arch_children_cnt      64037 non-null int64
14  arch_rm_cnt            64037 non-null int64
15  arch_destination_id    64037 non-null int64
16  arch_destination_type_id 64037 non-null int64
17  is_booking             64037 non-null int64
18  cnt                    64037 non-null int64
19  hotel_continent        64037 non-null int64
20  hotel_country          64037 non-null int64
21  hotel_market           64037 non-null int64
22  hotel_cluster          64037 non-null int64
23  stay_dur               64037 non-null float64
24  Cin_day               64037 non-null int64
25  Cin_month             64037 non-null int64
26  Cin_year              64037 non-null int64
dtypes: datetime64[ns](3), float64(2), int64(22)
memory usage: 13.7 MB
```

```
In [20]: # Count the bookings in each month
fig, ax = plt.subplots()
fig.set_size_inches(15, 8)
sns.countplot("Cin_month", data=df[df["is_booking"] == 1], order=list(range(1, 13)), ax=ax)
```

```
Out[20]:
```

FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be "data", and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
Out[20]: <AxesSubplot>
```

```
In [21]: # Count the bookings as per the day
fig, ax = plt.subplots()
fig.set_size_inches(15, 8)
sns.countplot("Cin_day", data=df[df["is_booking"] == 1], order=list(range(1, 32)), ax=ax)
```

```
Out[21]:
```

FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be "data", and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
Out[21]: <AxesSubplot>
```

```
In [22]: # Count the bookings as per the stay_duration
fig, ax = plt.subplots()
fig.set_size_inches(15, 8)
sns.countplot("stay_dur", data=df[df["is_booking"] == 1], ax=ax)
```

```
Out[22]:
```

FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be "data", and passing other arguments without an explicit keyword will result in an error or misinterpretation.

```
Out[22]: <AxesSubplot>
```

```
In [23]: # Check the percentage of NaN in dataset
#total = df.isnull().sum()
#percent = (df.isnull().sum()/df.shape[0])
#missing_data = pd.concat([total, percent], axis=1, keys=['Total', 'Percent'])
#missing_data
```

```
In [24]: # Fill nan with the day which has max occurrence
df["Cin_day"] = df["Cin_day"].fillna(26.0)
df["Cin_month"] = df["Cin_month"].fillna(8.0)
df["Cin_year"] = df["Cin_year"].fillna(2014.0)
df["stay_dur"] = df["stay_dur"].fillna(1.0)
```

```
# Fill average values in place for NaN, fill with mean
df["orig_destination_distance"].fillna(df["orig_destination_distance"].mean(), inplace=True)
```

```
In [25]: # Get list of categorical variable columns
catCols = ['site_name', 'posa_continent', 'user_location_country',
           'user_location_region', 'user_location_city',
           'channel', 'arch_destination_id', 'arch_children_cnt',
           'hotel_continent', 'hotel_country', 'hotel_market']
```

```
In [26]: # Function to convert to categorical
def df_cat(df):
    for col in catCols:
        df[col] = df[col].astype('category')
```

```
In [27]: # Converts a column to binary based on matching / not matching value
def cat_to_binary(row, col, val):
    if row[col] == val:
        return 1
    else:
        return 0
```

```
# Converts anything over a particular value to a certain value
def bin_val(row, col, val):
    if row[col] > val:
        return 1
    else:
        return 0
```

```
In [28]: # Data Transformations and new feature creation
df['site_name_2'] = df.apply(lambda row: cat_to_binary(row, 'site_name', 2), axis=1)
df['posa_continent_3'] = df.apply(lambda row: cat_to_binary(row, 'posa_continent', 3), axis=1)
df['user_location_country_66'] = df.apply(lambda row: cat_to_binary(row, 'user_location_country', 66), axis=1)
df['user_location_region'] = df.apply(lambda row: bin_val(row, 'user_location_region', 500), axis=1)
df['hotel_country'] = df.apply(lambda row: cat_to_binary(row, 'hotel_country', 50), axis=1)
```

```
In [29]: # Look at variables
df.info()
```

```
Out[29]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 64037 entries, 1 to 99999
Data columns (total 30 columns):
#   Column                Non-Null Count  Dtype
---  --
0   date_time              64037 non-null object
1   site_name              64037 non-null category
2   posa_continent         64037 non-null category
3   user_location_country  64037 non-null category
4   user_location_region   64037 non-null int64
5   user_location_city     64037 non-null category
6   orig_destination_distance 64037 non-null float64
7   is_mobile              64037 non-null int64
8   is_package             64037 non-null category
9   channel               64037 non-null object
10  arch_ci                64037 non-null object
11  arch_co                64037 non-null object
12  arch_children_cnt      64037 non-null int64
13  arch_rm_cnt            64037 non-null int64
14  arch_destination_id    64037 non-null int64
15  arch_destination_type_id 64037 non-null int64
16  is_booking             64037 non-null int64
17  cnt                    64037 non-null int64
18  hotel_continent        64037 non-null category
19  hotel_country          64037 non-null int64
20  hotel_market           64037 non-null category
21  hotel_cluster          64037 non-null int64
22  stay_dur               64037 non-null float64
23  Cin_day               64037 non-null int64
24  Cin_month             64037 non-null int64
25  Cin_year              64037 non-null int64
26  posa_continent_3      64037 non-null int64
27  user_location_country_66 64037 non-null int64
28  user_location_region   64037 non-null int64
29  hotel_country          64037 non-null int64
30  stay_dur               64037 non-null float64
dtypes: category(9), datetime64[ns](3), float64(2), int64(16)
memory usage: 12.1 MB
```

```
In [30]: df.columns
Out[30]: Index(['date_time', 'site_name', 'posa_continent', 'user_location_country',
           'user_location_region', 'user_location_city',
           'orig_destination_distance', 'is_mobile', 'is_package', 'channel',
           'arch_ci', 'arch_co', 'arch_children_cnt', 'arch_destination_type_id',
           'is_booking', 'cnt', 'hotel_continent', 'hotel_country', 'hotel_market',
           'hotel_cluster', 'stay_dur', 'Cin_day', 'Cin_month', 'Cin_year',
           'site_name_2', 'posa_continent_3', 'user_location_country_66'],
          dtype='object')
```

```
In [31]: # Load libraries
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
```

```
In [32]: final_df = df[['user_location_region', 'user_location_city',
                  'orig_destination_distance', 'is_mobile', 'is_package', 'channel',
                  'arch_rm_cnt', 'arch_children_cnt', 'arch_destination_type_id',
                  'arch_channel', 'arch_co', 'arch_children_cnt', 'arch_rm_cnt',
                  'is_booking', 'cnt', 'hotel_continent', 'hotel_country', 'hotel_market',
                  'hotel_cluster', 'stay_dur', 'Cin_day', 'Cin_month', 'Cin_year',
                  'site_name_2', 'posa_continent_3', 'user_location_country_66']].copy()
```

Split data for feature selection:

```
In [33]: # Set up features target sets
X = final_df.drop(["hotel_cluster"], axis=1)
y = final_df["hotel_cluster"]

# Encode the target variables
le = LabelEncoder()
y = le.fit_transform(y)

# Split the data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=7)

# Standardize features
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

Feature Selection

```
In [34]: # Feature selection using F-test
from sklearn.feature_selection import f_classif
f = f_classif(X_train, y_train[0])

In [35]: # Feature selection using mutual information
from sklearn.feature_selection import mutual_info_classif
mi = mutual_info_classif(X_train, y_train)
```

```
In [36]: # Feature selection using logistic regression
logreg = LogisticRegression(max_iter=500).fit(X_train, y_train)
```

```
In [37]: # Feature selection using light_gbm
from lightgbm import LGBMClassifier
lgbm = LGBMClassifier(
    objective = 'multiclass',
    metric = 'multi_logloss',
    importance_type = 'gain',
).fit(X_train, y_train)
```

```
Out[37]:
```

Traceback (most recent call last):

```
File ~/anaconda3/lib/python3.8/site-packages/lightgbm/_init_.py in <module>
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
```

```
Out[37]:
```

Library not loaded: /usr/local/opt/libomp.dylib
Referenced from: /Users/ykarki/opt/anaconda3/lib/python3.8/site-packages/lightgbm/lib_lightgbm.so
Reason: image not found

Store Results

```
In [38]: # Create DF to store feature ranking info
ranking = pd.DataFrame(index = range(X_train.shape[1]))

# Store ranking info for each feature from each method
ranking['f1'] = pd.Series(f, index = ranking.index).fillna(0).rank(ascending = False)
ranking['mi'] = pd.Series(mi, index = ranking.index).fillna(0).rank(ascending = False)
ranking['logreg'] = pd.Series(logreg.coef_.mean(axis = 0), index = ranking.index).rank(ascending = False)
ranking['lightgbm'] = pd.Series(lgbm.feature_importances_, index = ranking.index).rank(ascending = False)
ranking['nmcc'] = pd.Series((list(ranking[1, len(nmcc) + 1]) * (X_train.shape[1] - len(nmcc))
ranking = ranking.replace(to_replace = ranking[1, len(nmcc) + 1], value = X_train.shape[1])
ranking.to_csv('ranking.csv', index = False)
```

Compare

```
In [39]: # Run the rankings across methods
ranking['total'] = ranking.sum(axis=1)
# Sort by sum to get an overall idea of the most/least useful features
ranking.sort_values('total')
```

```
In [40]: # Get the highest ranked features
hiRank = ranking[ranking['total'] <= 110].sort_values('total')
```



```
# Create lists of models and accuracy scores
modelName = []
score = []
for key in models_dict:
    modelName += [key]
    model = models_dict[key]
    value = model.score(X_test, y_test)*100
    score += [round(value, 2)]

# Create DataFrame of results
d = {'Model': modelName, 'Accuracy': score}
results = pd.DataFrame(d).sort_values(by=['Accuracy'], ascending=False)
results
```

The Decision Tree model performed best with highly correlated variables, all others performed as good or better with highly ranked features, with the mix landing in-between.

In []:

```
# Set up dictionary for model results
models_dict = {}
```

The following were tuned using initial random sample of 100,000 records:

Decision Tree

```
In [ ]: # Load libraries
from sklearn.tree import DecisionTreeClassifier

# Create decision tree regressor object
decisiontree = DecisionTreeClassifier(random_state=42)

# Get a baseline model
baseline = decisiontree.fit(X_train, y_train)

# Create range of candidate penalty hyperparameter values
parameter_space = {
    'criterion': ['gini', 'entropy'],
    'splitter': ['best', 'random'],
    'max_depth': [None, 5, 8, 10],
    'min_samples_split': [2, 4, 8],
    'min_samples_leaf': [1, 2, 4],
    'max_leaf_nodes': [30, 100, 300, 1000, None],
    'class_weight': ['balanced', None],
    'max_features': ['sqrt', 'log2', None],
}

grid = GridSearchCV(decisiontree, parameter_space, verbose=2, n_jobs=-1, cv=5)

In [ ]: # Fit models
grid_result = grid.fit(X_train, y_train)

In [ ]: # Show best parameters
print('Best parameters found:\n', grid_result.best_params_, '\n')
# Get accuracy scores
baseScore = round(baseline.score(X_test, y_test)*100, 2)
score = round(grid_result.score(X_test, y_test)*100, 2)
print(f"Baseline Accuracy:\t{baseScore}")
print(f"Tuned Accuracy:\t{score}")
```

AdaBoost

```
In [ ]: # Load libraries
from sklearn.ensemble import AdaBoostClassifier

# Create classifier
adaboost = AdaBoostClassifier(random_state=42)

# Get a baseline model
baseline = adaboost.fit(X_train, y_train)

# Create range of candidate penalty hyperparameter values
parameter_space = {
    'n_estimators': [100, 300, 1000, 3000],
    'learning_rate': [0.1, 1, 10],
    'algorithm': ['SAMME', 'SAMME.R'],
}

grid = GridSearchCV(adaboost, parameter_space, n_jobs=-1, cv=5)

In [ ]: # Fit models
grid_result = grid.fit(X_train, y_train)

In [ ]: # Show best parameters
print('Best parameters found:\n', grid_result.best_params_, '\n')
# Get accuracy scores
baseScore = round(baseline.score(X_test, y_test)*100, 2)
score = round(grid_result.score(X_test, y_test)*100, 2)
print(f"Baseline Accuracy:\t{baseScore}")
print(f"Tuned Accuracy:\t{score}")
```

Random Forest

```
In [ ]: # Load libraries
from sklearn.ensemble import RandomForestClassifier

# Create classifier
rfc = RandomForestClassifier(random_state=42,
                             n_jobs=-1)

# Get a baseline model
baseline = rfc.fit(X_train, y_train)

# Create range of candidate penalty hyperparameter values
parameter_space = {
    'n_estimators': [30, 100, 300, 1000],
    'criterion': ['gini', 'entropy'],
    'max_features': ['sqrt', 'log2', None],
}

grid = GridSearchCV(rfc, parameter_space, n_jobs=-1, cv=5)

In [ ]: # Fit models
grid_result = grid.fit(X_train, y_train)

In [ ]: # Show best parameters
print('Best parameters found:\n', grid_result.best_params_, '\n')
# Get accuracy scores
baseScore = round(baseline.score(X_test, y_test)*100, 2)
score = round(grid_result.score(X_test, y_test)*100, 2)
print(f"Baseline Accuracy:\t{baseScore}")
print(f"Tuned Accuracy:\t{score}")
```

Regression

```
In [ ]: # Load libraries
from sklearn.linear_model import LogisticRegression

# Create logistic regression
logistic = LogisticRegression(random_state=42,
                              max_iter=100,
                              multi_class='auto',
                              n_jobs=-1)

# Get a baseline model
baseline = logistic.fit(X_train, y_train)

# Create range of candidate penalty hyperparameter values
parameter_space = {
    'penalty': ['l1', 'l2', 'elasticnet', 'none'],
    'tol': [1e-2, 1e-3, 1e-4],
    'C': np.logspace(0, 2, 3),
    'solver': ['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga'],
}

grid = GridSearchCV(logistic, parameter_space, n_jobs=-1, cv=5)

In [ ]: # Fit models
grid_result = grid.fit(X_train, y_train)

In [ ]: # Show best parameters
print('Best parameters found:\n', grid_result.best_params_, '\n')
# Get accuracy scores
baseScore = round(baseline.score(X_test, y_test)*100, 2)
score = round(grid_result.score(X_test, y_test)*100, 2)
print(f"Baseline Accuracy:\t{baseScore}")
print(f"Tuned Accuracy:\t{score}")
```

Support Vector Classifier (SVC)

```
In [ ]: # Load libraries
from sklearn.svm import SVC

# Create support vector classifier
svc = SVC(random_state=42,
          max_iter=10000)

# Get a baseline model
baseline = svc.fit(X_train, y_train)

# Create range of candidate penalty hyperparameter values
parameter_space = {
    'C': [1, 3, 10, 30, 100],
    'kernel': ['linear', 'poly', 'rbf', 'sigmoid'],
    'gamma': ['scale', 'auto'],
    'shrinking': (True, False),
    'tol': [1e-3, 1e-4, 1e-5],
    'decision_function_shape': ['ovo', 'ovr'],
    'break_ties': (True),
}

grid = GridSearchCV(svc, parameter_space, n_jobs=-1, cv=5)

In [ ]: # Fit models
grid_result = grid.fit(X_train, y_train)

In [ ]: # Show best parameters
print('Best parameters found:\n', grid_result.best_params_, '\n')
# Get accuracy scores
baseScore = round(baseline.score(X_test, y_test)*100, 2)
score = round(grid_result.score(X_test, y_test)*100, 2)
print(f"Baseline Accuracy:\t{baseScore}")
print(f"Tuned Accuracy:\t{score}")
```

Modeling with Tuned Parameters

```
# Set up dictionary for model results
models_dict = {}
```

Logistic Regression

```
In [ ]: # Load libraries
from sklearn.linear_model import LogisticRegression

# Create classification model
logistic = LogisticRegression(C=1.0,
                              penalty='l1',
                              solver='liblinear',
                              tol=0.001,
                              random_state=42,
                              n_jobs=-1)

# Fit model
models_dict['Logistic'] = logistic.fit(X_train, y_train)
```

Random Forest

```
In [ ]: # Load libraries
from sklearn.ensemble import RandomForestClassifier

# Create classification model
rfclassifier = RandomForestClassifier(criterion='entropy',
                                    max_features=None,
                                    n_estimators=1000,
                                    random_state=42,
                                    n_jobs=-1)

# Fit model
models_dict['RandomForest'] = rfclassifier.fit(X_train, y_train)
```

Decision Tree

```
In [ ]: # Load libraries
from sklearn.tree import DecisionTreeClassifier

# Create classification model
decisiontree = DecisionTreeClassifier(class_weight=None,
                                     criterion='entropy',
                                     max_depth=None,
                                     max_features=None,
                                     max_leaf_nodes=300,
                                     min_samples_leaf=1,
                                     min_samples_split=2,
                                     splitter='best',
                                     random_state=42,)

# Fit model
models_dict['DecisionTree'] = decisiontree.fit(X_train, y_train)
```

AdaBoost

```
In [ ]: # Load libraries
from sklearn.ensemble import AdaBoostClassifier

# Create classification model
adaboost = AdaBoostClassifier(algorithm='SAMME',
                              learning_rate=1,
                              n_estimators=1000,
                              random_state=42)

# Fit model
models_dict['AdaBoost'] = adaboost.fit(X_train, y_train)
```

Support Vector Classifier (SVC)

```
In [ ]: # Load libraries
from sklearn.svm import SVC

# Create classification model
svc = SVC(C=1,
          break_ties=True,
          decision_function_shape='ovr',
          gamma='scale',
          kernel='rbf',
          shrinking=True,
          tol=0.001,
          random_state=42)

# Fit model
models_dict['SVC'] = svc.fit(X_train, y_train)
```

Gaussian Naive Bayes Classifier

```
In [ ]: # Load libraries
from sklearn.naive_bayes import GaussianNB

# Create Gaussian naive Bayes object
nBayes = GaussianNB(var_smoothing=0.1)

# Fit model
models_dict['GaussianNB'] = nBayes.fit(X_train, y_train)
```

MLPClassifier

```
In [ ]: # Load libraries
from sklearn.neural_network import MLPClassifier

# Create classification model
mlp = MLPClassifier(tol=1e-5,
                    solve='adam',
                    learning_rate_init=0.01,
                    hidden_layer_sizes=(32, 68, 24),
                    epsilon=1e-3,
                    early_stopping=True,
                    alpha=1e-5,
                    activation='logistic')

# Fit model
models_dict['MLP'] = mlp.fit(X_train, y_train)

In [ ]:
```

Compare Accuracies

```
In [ ]: # Create lists of models and accuracy scores
modelName = []
score = []
for key in models_dict:
    modelName += [key]
    model = models_dict[key]
    value = model.score(X_test, y_test)*100
    score += [round(value, 2)]

# Create DataFrame of results
d = {'Model': modelName, 'Accuracy': score}
results = pd.DataFrame(d).sort_values(by=['Accuracy'], ascending=False)
results
```

In []: