**kaggle**    Search kaggle    Q    **Competitions    Datasets    Kernels    Discussion    Jobs**

# Expedia Hotel Recommendations

Which hotel type will an Expedia customer book?

$25,000  ·  1,974 teams  ·  a year ago

Overview    **Data**    Kernels    Discussion    Leaderboard    Rules                    Late Submission

---

**Competition Data**                                                                    Edit

📄 destinations.csv.gz         **train.csv.gz**  511.16 MB                    ⬇ Download

📄 sample_submission.cs...

📄 test.csv.gz

📄 train.csv.gz

---

**Data Description**

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

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

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

The train and test datasets are split based on time: training data from 2013 and 2014, while test data are from 2015. The public/private leaderboard data are split base on time as well. Training data includes all the users in the logs, including both click events and booking events. Test data only includes booking events.

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

Note that some srch_destination_id's in the train/test files don't exist in the destinations.csv file. This is because some hotels are new and don't have enough features in the latent space. Your algorithm should be able to handle this missing information.

# File descriptions

- **train.csv** - the training set
- **test.csv** - the test set
- **destinations.csv** - hotel search latent attributes
- **sample_submission.csv** - a sample submission file in the correct format

# Data fields

**train/test.csv**

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

**destinations.csv**

| Column name | Description | Data type |
|---|---|---|
| srch_destination_id | ID of the destination where the hotel search was performed | int |
| d1-d149 | latent description of search regions | double |

Our Team   Terms   Privacy   Contact/Support