

Overview

Search kaggle



Discussion Leaderboard

Datasets Kernels

Rules

Discussion



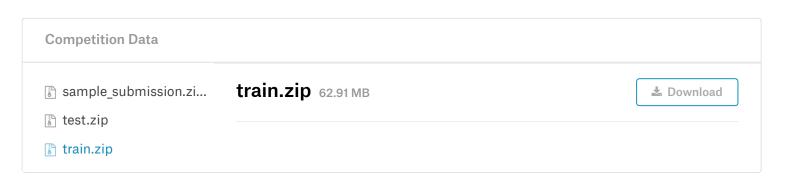
Late Submission

→ Playground Prediction Competition **New York City Taxi Trip Duration** \$30,000 Prize Money Share code and data to improve ride time predictions Kaggle · 1,257 teams · 2 days ago

You have accepted the rules for this competition.

Kernels

Data



## **Data Description**

The competition dataset is based on the 2016 NYC Yellow Cab trip record data made available in Big Query on Google Cloud Platform. The data was originally published by the NYC Taxi and Limousine Commission (TLC). The data was sampled and cleaned for the purposes of this playground competition. Based on individual trip attributes, participants should predict the duration of each trip in the test set.

## File descriptions

- train.csv the training set (contains 1458644 trip records)
- test.csv the testing set (contains 625134 trip records)
- sample submission.csv a sample submission file in the correct format

## **Data fields**

- id a unique identifier for each trip
- vendor id a code indicating the provider associated with the trip record
- pickup\_datetime date and time when the meter was engaged

- dropoff datetime date and time when the meter was disengaged
- passenger\_count the number of passengers in the vehicle (driver entered value)
- pickup\_longitude the longitude where the meter was engaged
- pickup\_latitude the latitude where the meter was engaged
- dropoff\_longitude the longitude where the meter was disengaged
- dropoff\_latitude the latitude where the meter was disengaged
- store\_and\_fwd\_flag This flag indicates whether the trip record was held in vehicle memory before sending to the vendor because the vehicle did not have a connection to the server Y=store and forward; N=not a store and forward trip
- trip\_duration duration of the trip in seconds

Disclaimer: The decision was made to not remove dropoff coordinates from the dataset order to provide an expanded set of variables to use in Kernels.

© 2017 Kaggle Inc

Our Team Terms Privacy Contact/Support





