



Featured Prediction Competition

# Zillow Prize: Zillow's Home Value Prediction (Zestimate)

\$1,200,000

Prize Money

Can you improve the algorithm that changed the world of real estate?



Zillow · 809 teams · 7 months to go (4 months to go until merger deadline)

Overview

**Data**

Kernels

Discussion

Leaderboard

More

My Submissions

Submit Predictions

## Training Data

properties\_2016.csv...

sample\_submission.cs...

**train\_2016.csv.zip**

zillow\_data\_dictiona...

**train\_2016.csv.zip** 634.09 KB

Download

## Data Introduction

In this competition, Zillow is asking you to predict the log-error between their Zestimate and the actual sale price, given all the features of a home. The log error is defined as

$$\text{logerror} = \log(\text{Zestimate}) - \log(\text{SalePrice})$$

and it is recorded in the transactions file **train.csv**. In this competition, you are going to predict the logerror for the months in Fall 2017. Since all the real estate transactions in the U.S. are publicly available, we will close the competition (no longer accepting submissions) before the evaluation period begins.

## Train/Test split

- You are provided with a full list of real estate properties in three counties (Los Angeles, Orange and Ventura, California) data in 2016.
- The train data has all the transactions before October 15, 2016, plus some of the transactions after October 15, 2016.
- The test data in the public leaderboard has the rest of the transactions between October 15 and December 31, 2016.
- The rest of the test data, which is used for calculating the private leaderboard, is **all** the properties in October 15, 2017, to December 15, 2017. This period is called the "sales tracking period", during which we will not be taking any submissions.
- You are asked to predict 6 time points for **all** properties: **October 2016 (201610), November 2016 (201611), December 2016 (201612), October 2017 (201710), November 2017 (201711), and December 2017 (201712)**.
- Not all the properties are sold in each time period. If a property was not sold in a certain time period, that particular row will be ignored when calculating your score.

- If a property is sold multiple times within 31 days, we take the first reasonable value as the ground truth. By "reasonable", we mean if the data seems wrong, we will take the transaction that has a value that makes more sense.

## File descriptions

- **properties\_2016.csv** - all the properties with their home features for 2016. Note: Some 2017 new properties don't have any data yet except for their parcelid's. Those data points should be populated when **properties\_2017.csv** is available.
- **properties\_2017.csv** - all the properties with their home features for 2017 (will be available on 10/2/2017)
- **train\_2016.csv** - the training set with transactions from 1/1/2016 to 12/31/2016
- **train\_2017.csv** - the training set with transactions from 1/1/2017 to 9/15/2017 (will be available on 10/2/2017)
- **sample\_submission.csv** - a sample submission file in the correct format

## Data fields

- Please refer to **zillow\_data\_dictionary.xlsx**